

Multi-level Searchable Symmetric Encryption

Sarah Louise Renwick
Information Security Group,
Royal Holloway, University of
London,
Egham, United Kingdom
sarahlouise.renwick.2012
@live.rhul.ac.uk

James Alderman
Information Security Group,
Royal Holloway, University of
London,
Egham, United Kingdom
james.alderman@rhul.ac.uk

Keith M. Martin
Information Security Group,
Royal Holloway, University of
London,
Egham, United Kingdom
keith.martin@rhul.ac.uk

ABSTRACT

Remote storage delivers a cost effective solution for data storage. If data is of a sensitive nature, it should be encrypted prior to outsourcing to ensure confidentiality; however, searching then becomes challenging. Searchable encryption is a well-studied solution to this problem. Many schemes only consider the scenario where users can search over the *entirety* of the encrypted data.

In practice, sensitive data is likely to be classified according to an access control policy and different users should have different access rights. It is unlikely that all users have unrestricted access to the entire data set. Current schemes that consider multi-level access to searchable encryption are predominantly based on asymmetric primitives. We investigate *symmetric* solutions to multi-level access in searchable encryption where users have different access privileges to portions of the encrypted data and are not permitted to search over, or learn information about, data for which they are not authorised.

Keywords

Searchable Symmetric Encryption, Multi-level, Access Control, Information Flow Policy

1. INTRODUCTION

Searchable encryption (SE) enables a user to search over encrypted data that has been outsourced to a remote server. In some schemes [4, 5, 9, 18, 19, 20], the data owner may authorise multiple users to make search queries — in such cases, a querier is either authorised to search over the entirety of the data or not at all, in which case (ideally) no information about the outsourced data should be revealed. In practice, the access control requirements of outsourced data sets are likely to be more fine-grained than this binary ‘all or nothing’ approach; hence existing schemes do not suffice.

We study the problem of enforcing a *multi-level access control policy (MLA)* in the context of searchable symmet-

ric encryption (SSE). As a notable example of this form of data classification, the UK government uses three levels of data classification: official, secret and top secret [16]. In our model, a user with ‘secret’ clearance should be unable to learn any information about data items classified as ‘top secret’, such as whether they contain searched keywords or not. This is an example of an information flow policy with a total order of security labels [2].

More precisely, consider a (possibly large) data set which is to be outsourced to an external storage provider, which could be outside of the data owner’s trusted zone. Although the provider has a business incentive to provide a storage and search service to the client (and to any other users authorised by the data owner), the provider may attempt to learn information about the sensitive data stored; in short, the storage provider may be *honest-but-curious*. Hence, the data must be encrypted prior to outsourcing, and the search procedure should not reveal unintended information to the storage provider or to other unauthorised entities. Each data item within the data set may be associated with some keywords, over which searches may be performed. Furthermore, each data item may differ in sensitivity and have different access control requirements. The data owner may authorise additional users to search the data set and, again, each user may have different access control clearance and therefore be able to access or search different sets of data items. Let us define a set of security labels \mathbb{L} , which forms a totally ordered set (\mathbb{L}, \leq) to reflect the inheritance of access rights. Each user u and data item d is assigned one of these labels, denoted $\lambda(u)$ and $\lambda(d)$ respectively. A user u may search a data item d if and only if $\lambda(u) \geq \lambda(d)$.

Public-key encryption (PKE), especially functional encryption, has previously been used to achieve MLA in SE [3, 11, 15, 21]. In general, PKE is computationally more intensive than symmetric key encryption (SKE), perhaps making SKE more suitable for practical systems. The enforcement of MLA policies in *symmetric* SE has, up to now, remained relatively unexplored. Kissel et al. [14] presented a SKE-based scheme in which users are divided into groups that each have a specified dictionary of keywords they may search over. These groups are arranged hierarchically so that each group may also search for all keywords in dictionaries assigned to groups at lower levels in the hierarchy. Although this scheme presents a form of hierarchical access in SSE, users may still search over the entire data set. In most access control scenarios, we are concerned with protecting a data item (i.e. the complete content of a data item), not just a single keyword describing the data item. Furthermore,

it may be difficult to correctly administer an access control policy expressed only in terms of authorised keywords; data items may gain their classification level due to semantic meaning regarding their contents (for example, the subject to which they pertain), which may not trivially be captured through the associated keywords. For example, consider two data items containing information about company spending: one providing a public report of company-wide spending, whilst the other pertains specifically to the research department. Whilst both items may be labelled by a keyword such as ‘finance’, detailed knowledge of research spending may be deemed more sensitive than a generalised report. Simply authorising users to search for keywords, such as ‘finance’, does not suffice in this instance as not all users that can search the public report should also be able to view the specific report. The access control policy in this case must be managed carefully — perhaps additional, more granular, keywords must be defined e.g. ‘finance-public’ (leading to an increase in the size of the searchable encryption index and a subsequent loss of efficiency) or a (less efficient) SE scheme that supports ‘conjunctive keyword-only access control’ would be required such that one can be authorised to search for (‘finance’ AND ‘public’) and only data items with *both* keywords would be returned. In this work, we consider the problem of fine-grained classification of data items *directly* and gain a more efficient solution.

In this work, we consider Multi-level Searchable Symmetric Encryption (MLSSE). We begin in Section 2 by reviewing background material, before defining our system and security models in Sections 3.1 and 3.2. In Section 3.3, we introduce our instantiation based on the constructions of [9, 13], and then show, in Section 3.5, how to extend our construction to support a dynamic data set using techniques from [13]. Section 3.6 discusses the efficiency of our scheme. The full security proofs of our constructions are omitted but are available in the full version of our paper [1].

2. BACKGROUND

We aim to enforce *information flow policies* within searchable encryption, which encompass a wide range of access control policies that are of practical interest, including the Bell-LaPadula model, temporal, role-based and attribute-based access control [8].

Definition 1. An *information flow policy* is a tuple $\mathcal{P} = ((\mathbb{L}, \leq), \mathcal{U}, \mathcal{D}, \lambda)$, where (\mathbb{L}, \leq) is a partially ordered set (poset)¹ of security labels, \mathcal{U} is a set of users, \mathcal{D} is a set of objects (data items), and $\lambda : \mathcal{U} \cup \mathcal{D} \rightarrow \mathbb{L}$ is a function mapping users and objects to security labels in \mathbb{L} . We say that $u \in \mathcal{U}$ is *authorised* to read (search) an object $d \in \mathcal{D}$ if $\lambda(d) \leq \lambda(u)$.

In this paper, we will focus on the case where (\mathbb{L}, \leq) is a *total order* (chain) giving a simple hierarchy of security levels and, without loss of generality, we assume that each user and object is assigned to at most one security label. Given a set X , we denote the power set of X , comprising all combinations of elements in X , by 2^X . Throughout this paper we refer to ‘security levels’ and ‘security labels’ as *access levels*.

¹A poset is a set of labels L and a binary order relation \leq on L such that for all x, y and $z \in L$, $x \leq x$ (reflexivity), if $x \leq y$ and $y \leq x$ then $x = y$ (antisymmetry), and if $x \leq y$ and $y \leq z$ then $x \leq z$ (transitivity). If $x \leq y$ then we may write $y \geq x$.

Definition 2. A *Multi-User Searchable Symmetric Encryption (MSSE)* scheme is a set of six polynomial time algorithms defined as follows:

- $K_O \xleftarrow{\$} \text{MSSE.KeyGen}(1^k)$: A probabilistic algorithm run by the data owner that takes a security parameter $k \in \mathbb{N}$ and outputs a secret key K_O .
- $(\mathcal{I}_D, st_O, st_S) \xleftarrow{\$} \text{MSSE.BuildIndex}(K_O, \mathcal{D}, \mathcal{G})$: A probabilistic algorithm run by the data owner that takes a set of data items \mathcal{D} , a set of authorized users \mathcal{G} and the secret key K_O . It outputs an index \mathcal{I}_D , and server and owner states st_S and st_O .
- $K_u \xleftarrow{\$} \text{MSSE.AddUser}(u, K_O, st_O)$: A probabilistic algorithm run by the data owner that takes the identity, u , of a user to be enrolled in the system along with the owner’s secret key and state. It outputs a secret key for the new user K_u .
- $T_\omega \leftarrow \text{MSSE.Query}(\omega, K_u)^2$: A deterministic algorithm run by a user that takes a keyword ω and the user’s secret key, and outputs a search token.
- $R_\omega \leftarrow \text{MSSE.Search}(T_\omega, \mathcal{I}_D, st_S)$: A deterministic algorithm run by the server that takes as input a search token, an encrypted index and the server state, and outputs a set R_ω of identifiers of data items containing ω .
- $(st_O, st_S) \xleftarrow{\$} \text{MSSE.Revoke}(u, K_O, st_O)$: A probabilistic algorithm run by the data owner that takes a user identity of a user to be revoked along with the data owner’s secret key and state. It outputs new server and owner states.

For a data set \mathcal{D} and keyword $\omega \in \Delta$ (where Δ is a dictionary of possible keywords), let us denote by \mathcal{D}_ω the expected results of searching for ω in \mathcal{D} (in the plain); informally we say that an MSSE scheme is correct if it also produces the output \mathcal{D}_ω . More formally, a MSSE scheme MSSE is correct if for all $k \in \mathbb{N}$, for all K_O output by $\text{MSSE.KeyGen}(1^k)$, for all $\mathcal{D} \in 2^\Delta$, for all $\mathcal{G} \in 2^\mathcal{U}$, for all $(\mathcal{I}_D, st_O, st_S)$ output by $\text{MSSE.BuildIndex}(K_O, \mathcal{G}, \mathcal{D})$, for all ω in Δ : $\text{Search}(\text{MSSE.Query}(K_u, \omega), \mathcal{I}_D, st_S) = \mathcal{D}_\omega$.

Definition 3. A *Broadcast encryption (BE)* scheme is a set of four polynomial time algorithms as follows, where \mathcal{U} is the user space of all possible user identities:

- $(PP, K_{BE}) \xleftarrow{\$} \text{BE.Keygen}(1^k)$: A probabilistic algorithm that takes a security parameter k outputs public parameters PP and a master secret key K_{BE} .
- $C \xleftarrow{\$} \text{BE.Enc}(M, \mathcal{G})$: A probabilistic algorithm that takes a plaintext M , a set of users $\mathcal{G} \in \mathcal{U}$ authorized to decrypt and produces a ciphertext C .
- $K_u \xleftarrow{\$} \text{BE.Add}(K_{BE}, u)$: A probabilistic algorithm that takes as input the master secret key K_{BE} and a user identifier $u \in \mathcal{U}$, and outputs a user key K_u .

²This algorithm is sometimes referred to as MSSE.Trapdoor in the literature, however to maintain consistent notation throughout this paper we refer to it as MSSE.Query

- $(M \text{ or } \perp) \leftarrow \text{BE.Dec}(C, K_u)$: A deterministic algorithm that takes a ciphertext C and a secret key K_u and outputs either a plaintext M or a failure symbol \perp .

BE is correct if $\forall k \in \mathbb{N}$, for all PP and K_{BE} output by $\text{BE.KeyGen}(1^k, m)$, for all M in the plaintext space, all sets of users $\mathcal{G} \in \mathcal{U}$, every K_u output by $\text{BE.Add}(u, K_{\text{BE}})$ and all C output by $\text{BE.Enc}(M, \mathcal{G})$ where $u \in \mathcal{G}$ we have: $M \leftarrow \text{BE.Dec}(C, K_u)$.

3. MULTI-LEVEL ACCESS IN SEARCHABLE SYMMETRIC ENCRYPTION

A MLSSE scheme permits searching over encrypted data in the symmetric key setting for multiple users that have varying access rights to the set of data items. The access levels are hierarchical (totally ordered), meaning a user may search all data items at their own access level as well as all data items that are classified at lower access levels.

3.1 System Model

Consider a *data owner* O , a *server* S , and a set of m data users $\mathcal{U} = \{u_1, \dots, u_m\}$. The data owner possesses a set of data items $\mathcal{D} = \{d_1, \dots, d_n\}$ which they wish to encrypt and outsource to S whilst authorising other users to search over some data items within \mathcal{D} . Each data item $d_i \in \mathcal{D}$ is associated with an identifier id_{d_i} .

To enable searching over the encrypted data, O must upload some encrypted metadata to the server. It first defines a dictionary of keywords, denoted $\Delta = \{\omega_1, \dots, \omega_{|\Delta|}\}$, and assigns a set $\delta_{d_i} \subseteq \Delta$ of keywords to each data item $d_i \in \mathcal{D}$. We refer to the set of keywords for all data items as $\delta_{\mathcal{D}} = (\delta_{d_1}, \dots, \delta_{d_n})$. The data owner then produces an encrypted *index* $\mathcal{I}_{\mathcal{D}}$ based on $\delta_{\mathcal{D}}$, over which searches will be performed.

O also defines an information flow policy \mathcal{P} with a labelling function λ mapping each user $u_i \in \mathcal{U}$ and data item $d_j \in \mathcal{D}$ to an access level, denoted $\lambda(u_i)$ and $\lambda(d_j)$ respectively, in the totally ordered set $\mathbb{L} = \{a_1, \dots, a_l\}$. Access control in our model is enforced at data item level — users are restricted in the data items that they may search, not the keywords they may search for [14]. A user with access level $\lambda(u_i)$ is authorised to search a data item with classification $\lambda(d_j)$ if and only if $\lambda(d_j) \leq \lambda(u_i)$. To search for a keyword $\omega \in \Delta$, a user u_i (with access level $\lambda(u_i)$) generates a search query $T_{\omega, \lambda(u_i)}$. Let \mathcal{D}_{ω} be the set of identifiers of all data items assigned the keyword ω , and denote by $\mathcal{D}_{\omega, \lambda(u_i)} \subseteq \mathcal{D}_{\omega}$ the search results that user u_i is authorised to view; in other words, the set of identifiers of all data items id_{d_j} assigned ω where $\lambda(d_j) \leq \lambda(u_i)$.

To add and revoke users, we use *broadcast encryption* (BE) (Definition 3) as per [9]; a user may only produce a valid search query if they are authorized in the BE scheme.

To ease notation, we define the tuple $d_i^{aug} = (d_i, id_i, \delta_{d_i}, \lambda(d_i))$ to completely describe a data item $d_i \in \mathcal{D}$ (being the data itself, the identifier, the associated keywords and the security classification). We denote the information regarding all data items by $\mathcal{D}^{aug} = \{d_1^{aug}, \dots, d_n^{aug}\}$.

We present a *structure only* MLSSE system — we only consider the data structure (index) and do not encrypt the data items themselves; data items may be encrypted separately and retrieved based on the search results, which comprise a set of data item identifiers that fulfil the query. We

permit data items to be of any format and the sets of keywords can be arbitrarily chosen from the dictionary — they may not necessarily correspond to the actual content of the data, but could be descriptive attributes of the data item. This may help minimise the risk of a statistical attack on the index as the frequency of a certain word in a document is not necessarily reflected in the set of keywords chosen to index the data item.

Definition 4. A *Multi-level Searchable Symmetric Encryption Scheme (MLSSE)* scheme consists of six algorithms defined as follows:

- $(K_O, k_S, PP) \stackrel{\$}{\leftarrow} \text{KeyGen}(1^k, S, \mathcal{P})$: A probabilistic algorithm run by the data owner O that takes the security parameter k , policy \mathcal{P} and the server identity S , and outputs O 's secret key K_O , a server key k_S and public parameters PP .
- $\mathcal{I}_{\mathcal{D}} \stackrel{\$}{\leftarrow} \text{BuildIndex}(\mathcal{D}^{aug}, K_O, PP)$: A probabilistic algorithm run by O . It takes the description of the data set \mathcal{D}^{aug} and O 's secret key, and outputs the index $\mathcal{I}_{\mathcal{D}}$.
- $(K_u, PP) \stackrel{\$}{\leftarrow} \text{AddUser}(u, \lambda(u), K_O, PP)$: A probabilistic algorithm run by O to enrol a new user into the system. It takes the new user's identity u and access level $\lambda(u)$, and O 's key, and outputs a secret key for the new user.
- $T_{\omega, \lambda(u)} \leftarrow \text{Query}(\omega, K_u)$: A deterministic algorithm run by a user with access level $\lambda(u)$ to generate a search query. It takes as input a keyword $\omega \in \Delta$ and the user's secret key and outputs a search query $T_{\omega, \lambda(u)}$.
- $\mathcal{R}_{\omega, \lambda(u)} \leftarrow \text{Search}(T_{\omega, \lambda(u)}, \mathcal{I}_{\mathcal{D}}, k_S)$: A deterministic algorithm run by S to search the index for data items containing a keyword ω . It takes a search query and the index, and returns the search results $\mathcal{R}_{\omega, \lambda(u)}$, comprising either a set $\mathcal{D}_{\omega, \lambda(u)}$ of identifiers of data items d_j containing ω such that for all $\lambda(d_j) \leq \lambda(u)$ (where $\lambda(u)$ is the access level of the user that submitted the search query), or a failure symbol \perp .
- $(K_O, PP) \stackrel{\$}{\leftarrow} \text{RevokeUser}(u, K_O, PP)$: A probabilistic algorithm run by O to revoke a user from the system. It takes the user's id, the data owner's and server's secret keys, and outputs updated owner and server keys.

An MLSSE scheme is correct if for all $k \in \mathbb{N}$, for all K_O, k_S output by $\text{KeyGen}(1^k, S, \mathcal{P})$, for all \mathcal{D}^{aug} , for all $\mathcal{I}_{\mathcal{D}}$ output by $\text{BuildIndex}(\mathcal{D}^{aug}, K_O, PP)$, for all $\omega \in \Delta$, for all $u \in \mathcal{U}$, for all K_u output by $\text{AddUser}(u, \lambda(u), K_O, PP)$, $\text{Search}(\text{Query}(\omega, K_u), \mathcal{I}_{\mathcal{D}}, k_S) = \mathcal{D}_{\omega, \lambda(u)}$.

3.2 Security model

A secure MLSSE scheme would, ideally, reveal no information regarding the data set \mathcal{D} to the server (i.e. a curious server cannot learn information about the data it stores) and reveal no information to users regarding data items that they are not authorised to search. However, most SSE schemes leak additional information to gain efficiency. For example, the search results $\{R_{\omega_1, a}, \dots, R_{\omega_p, a}\}$ for a set of queries $\{T_{\omega_1}, \dots, T_{\omega_p}\}$ could be revealed. This is referred to as the *access pattern* (Definition 5) and defines the link between

a search query and the search results it produces; it may be thought of as a database where each row stores a search query and a corresponding set of identifiers of data items that satisfies the search query.

Most efficient SSE schemes also leak the *search pattern* (Definition 6), which reveals the set of search queries made to the server. In most single-user SSE schemes [6, 7, 9, 10, 12, 13], search queries are formed deterministically; the server can therefore ascertain whether a search query has been made previously.

Definition 5. For a sequence of q search queries $\Omega = \{T_{\omega_1, a_1}, \dots, T_{\omega_q, a_q}\}$ where for $1 \leq i, j \leq q$: ω_i and ω_j or $\lambda(u_i)$ and $\lambda(u_j)$ are not necessarily distinct for $i \neq j$, the *access pattern* is defined as:

$$AP(\mathcal{I}_{\mathcal{D}}, \Omega) = \{(T_{\omega_1, a_1}, R_{\omega_1, a_1}), \dots, (T_{\omega_q, a_q}, R_{\omega_q, a_q})\}.$$

Definition 6. For a sequence of q search queries $\Omega = \{T_{\omega_1, a_1}, \dots, T_{\omega_q, a_q}\}$ where for $1 \leq i, j \leq q$: ω_i and ω_j or $\lambda(u_i)$ and $\lambda(u_j)$ are not necessarily distinct for $i \neq j$, the *search pattern* is defined as a $q \times q$ symmetric binary matrix $SP(\mathcal{I}_{\mathcal{D}}, \Omega)$ such that for $1 \leq i, j \leq q$:

$$SP(\mathcal{I}_{\mathcal{D}}, \Omega)_{i,j} = 1 \iff T_{\omega_i, a_i} = T_{\omega_j, a_j}.$$

Intuitively, the *search pattern* reveals when the i th and j th queries are the same, which happens when queries are issued for the same keyword by users with the same access level.

Definition 7. For an index $\mathcal{I}_{\mathcal{D}}$ we define the *setup leakage* $\mathcal{L}_{Setup}(\mathcal{I}_{\mathcal{D}})$ to be all the information that is leaked by the index $\mathcal{I}_{\mathcal{D}}$.

Definition 8. For an index $\mathcal{I}_{\mathcal{D}}$ and set of q search queries $\Omega = (T_{\omega_1}, \dots, T_{\omega_q})$ we define the *query leakage* $\mathcal{L}_{Query}(\mathcal{I}_{\mathcal{D}}, \Omega)$ to be all the information leaked by evaluating the queries in Ω on the index $\mathcal{I}_{\mathcal{D}}$.

We now formalise the notions of security we require in MLSSE. We use cryptographic games to formalize our notions of security. For each game, a challenger \mathcal{C} instantiates a probabilistic polynomial time (PPT) adversary \mathcal{A} whose inputs are chosen to reflect the information available to a realistic adversary. Our notion of adaptive security is based on that of IND-CKA2 presented in [9]. In the following we represent the dictionary of keywords as Δ , λ defines the mapping function as described in Section 3.

3.2.1 Multi-level Access

Our first security notion, in Figure 1, is that of *multi-level access* which requires that a user, u , cannot receive search results or learn information relating to data items d_i such that $\lambda(u) < \lambda(d_i)$. More specifically, a server colluding with several users cannot learn anything about the index beyond the specified leakage according to the corrupt users' access rights.

We define a *maximal query leakage with access level* λ_{max} on $\mathcal{I}_{\mathcal{D}}$ to be $\mathcal{L}_{Query}(\mathcal{I}_{\mathcal{D}}, \{T_{\omega_i, \lambda_{max}}\}_{\omega_i \in \Delta})$ — this is the leakage resulting from every possible keyword search with the maximal access level available, in Game 1 we denote this as $\mathcal{L}^{max}(\mathcal{I}_{\mathcal{D}})$.

The challenger sets up the system, including instantiating several global variables (which the challenger can use in the main game and in oracle functions, but which the adversary

cannot see): L is a list of users that have been corrupted, λ^{max} is the maximal access level of any corrupted user, and $chall$ is a Boolean flag to show whether the challenge parameters have been generated yet. The adversary is given the security parameter, access control policy, server key and the public parameters, as well as providing access to the following oracles.

The ADDUSER oracle allows the adversary to enrol a user into the system, and the adversary corrupts this user by receiving the user key. If the challenge has not yet been generated, then the challenger adds the requested user to the list L of corrupted users, checks if the maximal access level of corrupted users needs updating, and runs the AddUser algorithm. Otherwise, if the challenge has been generated, the above procedure is carried out only if the maximal query leakage for the new user's access level is equal on both challenge data sets — that is, providing the user key for the queried user cannot allow the adversary to trivially distinguish the two data sets.

The REVOKEUSER oracle first checks that the requested user has indeed been added previously. If so, it removes the user identity from L and checks whether the maximal access level needs changing. It returns the server key resulting from running the RevokeUser algorithm.

The BUILDINDEX oracle simply runs BuildIndex and returns the output to the adversary.

After a polynomial number of queries, the adversary outputs two data sets which must have identical maximal query leakages for the maximal access level of any corrupted user. The adversary cannot choose data sets where a user that it has corrupted could make any query that legitimately distinguishes the data sets since this would count as a trivial win. Whilst this may appear to be a strong assumption, we believe it to be the minimal assumption necessary to avoid trivial wins in the multi-user setting. The main issue is that in the multi-user setting it is necessary to consider the server colluding with a set of users (but not the data owner); as such, the adversary is able to perform the roles of the server and of an authorised user, and therefore may produce arbitrary search queries and perform searches themselves. Thus, the challenger in the game is unable to monitor which searches have been performed and hence cannot determine whether the query leakages of the *actual* queries on both data sets are equal, and instead must rely on the stronger assumption that no possible authorised search query can distinguish the data sets. Note that Van Rompay et al. [17] deal with the multi-user case without this assumption since they deal with single word indexes and have a proxy through which all queries are made.

The challenger sets the challenge flag to true and chooses a random bit b which determines the data set used to form an index. The adversary is given the index and oracle access as described in Game 1 and must determine which data set was used.

Definition 9. (Multi-level Access) Let \mathcal{ML} be a multi-level searchable symmetric encryption scheme where $k \in \mathbb{N}$ is the security parameter, \mathcal{P} is an information flow policy, S is the identity of the server and \mathcal{A} a PPT adversary. The advantage of \mathcal{A} is:

$$Adv_{\mathcal{A}}^{MLA}(\mathcal{ML}, 1^k, \mathcal{P}) = |\Pr[\mathbf{Exp}_{\mathcal{A}}^{MLA}[\mathcal{ML}, 1^k, S, \mathcal{P}] = 1] - \frac{1}{2}|.$$

We say that \mathcal{ML} is $(\mathcal{L}_{Setup}, \mathcal{L}_{Query})$ -secure against adaptive

Game 1 $\text{Exp}_{\mathcal{A}}^{MLA}[X]$:

```
1:  $\lambda^{max} \leftarrow \perp$ 
2:  $\text{chall} \leftarrow \text{false}$ 
3:  $(K_O, K_S, PP) \xleftarrow{\$} \text{KeyGen}(1^\kappa, \mathcal{U}, \mathcal{P})$ 
4:  $(\mathcal{D}_0^{aug}, \mathcal{D}_1^{aug}, st) \xleftarrow{\$} \mathcal{A}^O(X, K_S, PP)$ 
5:  $\mathcal{I}_{\mathcal{D}_0} \xleftarrow{\$} \text{BuildIndex}(\mathcal{D}_0^{aug}, K_O, PP)$ 
6:  $\mathcal{I}_{\mathcal{D}_1} \xleftarrow{\$} \text{BuildIndex}(\mathcal{D}_1^{aug}, K_O, PP)$ 
7: if  $\mathcal{L}_{Search}^{\lambda^{max}}(\mathcal{I}_{\mathcal{D}_0}) \neq \mathcal{L}_{Search}^{\lambda^{max}}(\mathcal{I}_{\mathcal{D}_1})$  then
8:   return 0
9:  $\text{chall} \leftarrow \text{true}$ 
10:  $b \xleftarrow{\$} \{0, 1\}$ 
11:  $b' \xleftarrow{\$} \mathcal{A}^O(\mathcal{I}_{\mathcal{D}_b}, st)$ 
12: if  $b' = b$  then return 1
13: else return 0
```

Oracle 1 $\text{AddUser}(u, \lambda(u), K_O, PP)$

```
1: if  $\text{chall} = \text{false}$  then
2:   if  $\lambda(u) > \lambda^{max}$  then
3:      $\lambda^{max} \leftarrow \lambda(u)$ 
4:     return  $\text{AddUser}(u, \lambda(u), K_O, PP)$ 
5:   else
6:     if  $\lambda(u) > \lambda^{max}$  then
7:       return  $\perp$ 
8:   else
9:     return  $\text{AddUser}(u, \lambda(u), K_O, PP)$ 
```

chosen keyword attacks in the sense of Game 1 if for all \mathcal{A} , all $k \in \mathbb{N}$, all S and all \mathcal{P} , $\text{Adv}_{\mathcal{A}}^{MLA}(\mathcal{ML}, 1^\kappa, S, \mathcal{P}) \leq \text{negl}(k)$ for a negligible function negl .

3.2.2 Revocation Security

In MLSSE, as with other multi-user SSE schemes, we need to consider user *revocation* to remove a user's ability to submit valid search queries to the server, and hence receive search results. We capture this in Game 2. The adversary is given the public parameters and selects a data set (along with associated access levels, keywords and identifiers). The challenger then creates the index. The adversary is given access to a set of oracles that perform the $\text{AddUser}(\cdot, \lambda(\cdot), K_O, PP)$, $\text{Search}(\cdot, \mathcal{I}_{\mathcal{D}}, k_S)$ and $\text{RevokeUser}(\cdot, K_O, PP)$ functions, where the parameters represented by \cdot are provided by the adversary, and the adversary is given the resulting user keys and search results. Once the adversary has completed his queries, the challenger revokes all users that were queried to the AddUser oracle but were not subsequently queried to the RevokeUser oracle (i.e. all users for which the adversary holds a valid user key). The adversary must then produce a search query T which, when used as input to the Search algorithm, does not produce \perp i.e. the adversary must produce a valid search query even though it does not hold a non-revoked key.

Definition 10. (Revocation) Let \mathcal{ML} be a multi-level searchable symmetric encryption scheme where $k \in \mathbb{N}$ is the security parameter, S the server identity, \mathcal{P} is an information flow policy and \mathcal{A} a PPT adversary. We define the advantage of \mathcal{A} in Game 2 as:

$$\text{Adv}_{\mathcal{A}}^{\text{Revoke}}(\mathcal{ML}, 1^\kappa, S, \mathcal{P}) = |\mathbb{P}[\text{Exp}_{\mathcal{A}}^{\text{Revoke}}[\mathcal{ML}, 1^\kappa, S, \mathcal{P}] = 1] - \frac{1}{2}|.$$

Game 2 $\text{Exp}_{\mathcal{A}}^{\text{Revoke}}[X]$:

```
1:  $\mathcal{G} \leftarrow \emptyset$ 
2:  $(K_O, K_S, PP) \xleftarrow{\$} \text{KeyGen}(1^\kappa, \mathcal{U}, \mathcal{P})$ 
3:  $(\mathcal{D}^{aug}, st) \leftarrow \mathcal{A}^O(X, PP)$ 
4:  $\mathcal{I}_{\mathcal{D}} \xleftarrow{\$} \text{BuildIndex}(\mathcal{D}^{aug}, K_O, PP)$ 
5:  $st \xleftarrow{\$} \mathcal{A}^O(st)$ 
6: for  $u \in \mathcal{G}$  do
7:    $(K_O, PP) \xleftarrow{\$} \text{RevokeUser}(u, K_O, PP)$ 
8:  $T_\omega \xleftarrow{\$} \mathcal{A}^O(st, PP)$ 
9:  $\mathcal{R} \leftarrow \text{Search}(T_\omega, \mathcal{I}_{\mathcal{D}}, K_S)$ 
10: if  $\mathcal{R} \neq \perp$  then
11:   return 1
12: else
13:   return 0
```

Oracle 2 $\text{AddUser}(u, \lambda(u), K_O, PP)$

```
1: if  $u \in \mathcal{G}$  then
2:   return  $\perp$ 
3: else
4:    $\mathcal{G} \leftarrow \mathcal{G} \cup u$ 
5:   return  $\text{AddUser}(u, \lambda(u), K_O, PP)$ 
```

Oracle 3 $\text{RevokeUser}(u, K_O, PP)$

```
1: if  $u \in \mathcal{G}$  then
2:    $\mathcal{G} \leftarrow \mathcal{G} \setminus u$ 
3:   return  $\text{RevokeUser}(u, K_O, PP)$ 
4: else
5:   if  $u \notin \mathcal{G}$  then
6:     return  $\perp$ 
```

We say that \mathcal{ML} achieves revocation if for all \mathcal{A} , all $k \in \mathbb{N}$, all S and all \mathcal{P} ,

$$\text{Adv}_{\mathcal{A}}^{\text{Revoke}}(\mathcal{ML}, 1^\kappa, S, \mathcal{P}) \leq \text{negl}(k).$$

3.3 Construction

Our construction MLSSE is an adaptation of the scheme of Kamara et al. [13], which is based on the construction of the influential inverted index scheme SSE-1 by Curtmola et al. [9].

Informally, MLSSE scheme uses an array \mathbb{A} of linked lists, along with a look-up table \mathbb{T} to index the encrypted data. This produces a sequential search that lends itself well to the hierarchical access rights on the data items that we require. For each keyword $\omega_i \in \Delta$, we define a list L_{ω_i} which stores the identifiers for all data items containing that keyword and is ordered according to the access level of the data items — data items with the highest classification are placed at the beginning of the list, and those with the lowest classification at the end. Each list L_{ω_i} is encrypted and stored in \mathbb{A} as a linked list. During the search phase the look-up table \mathbb{T} is used to point the server to the correct node in the array depending on the information in the search query i.e. which keyword was searched for and what access rights the user that submitted the search query has. This node is decrypted using information in the search query and the node itself, revealing the address of the next node in the linked list. The server may continue to decrypt all other relevant nodes in the linked list, obtaining the set of search results relevant to the user's searched keyword and access level.

The key difference between our scheme and that of [13] is that, rather than pointing to the *beginning* of each linked list, the entry in \mathbb{T} will point to the appropriate position within the linked list according to the access rights of the querier (recall that the list is ordered by access levels). Since it is not possible to move backwards through the encrypted lists, the only search results available are those contained beyond this point in this list — that is, identifiers for those documents containing the keyword and whose classification is at most that of the querier, as required by the information flow policy.

Let BE be an IND-CPA secure broadcast encryption scheme. We define the following pseudorandom functions (PRFs):

$$F : \{0, 1\}^k \times \{0, 1\}^* \rightarrow \{0, 1\}^k,$$

$$G : \{0, 1\}^k \times \{0, 1\}^* \rightarrow \{0, 1\}^*,$$

$$P : \{0, 1\}^k \times \{0, 1\}^* \rightarrow \{0, 1\}^k,$$

$$H : \{0, 1\}^* \times \{0, 1\}^k \rightarrow \{0, 1\}^*,$$

and a pseudorandom permutation (PRP):

$$\phi : \{0, 1\}^k \times \{0, 1\}^* \times \{0, 1\}^k \times \{0, 1\}^k \rightarrow \{0, 1\}^k \times \{0, 1\}^* \times \{0, 1\}^k,$$

\mathbb{A} is a $|\Delta| \times |\mathbb{L}|$ array and \mathbb{T} is a dictionary of size $|\Delta| \cdot |\mathbb{L}|$. We denote the address of a node N in \mathbb{A} as $\text{addr}_{\mathbb{A}}(N)$.

Let λ map users and data items to their relevant access levels as described in Section 3.1. We define a function γ which outputs three ordered lists \mathbb{L}_{ω_i} , \mathbb{X}_{ω_i} and \mathbb{N}_{ω_i} given the set of identifiers \mathcal{D}^{aug} and the array \mathbb{A} . We refer to the n^{th} item in a list \mathbb{L}_{ω_i} as $\mathbb{L}_i[n]$. The list \mathbb{L}_{ω_i} contains identifiers of data items in \mathcal{D}_{ω_i} ordered from the identifiers with the highest to the lowest access levels, the list \mathbb{N}_{ω_i} contains the addresses of $|\mathbb{L}_{\omega_i}|$ nodes chosen randomly from \mathbb{A} and the list \mathbb{X}_{ω_i} contains the indices of the identifiers in \mathbb{L}_{ω_i} where each access level starts i.e. if we have an ordered list of identifiers $\mathbb{L}_{\omega_i} = (id_1, id_2, id_3, id_4, id_5)$ where:

$$a_1 = \lambda(id_1) = \lambda(id_2) = \lambda(id_3) > \lambda(id_4) = \lambda(id_5) = a_3.$$

We have that $\mathbb{X}_{\omega_i}[3] = 4$, which says that the list of nodes with access level at most a_3 starts at the fourth entry in \mathbb{L}_{ω_i} . There is an entry per each access level in \mathbb{X}_{ω_i} , even if two access levels have the same starting point in \mathbb{L}_{ω_i} ; from the example above we can see that $\mathbb{X}_{\omega_i}[2] = \mathbb{X}_{\omega_i}[3] = 4$. If an access level is not authorised to view any data items in \mathcal{D}_{ω_i} then the entry corresponding to that access level (as well as the entries corresponding to all access levels below it) in \mathbb{X}_{ω_i} is set to \perp . An identifier of a data item $d_i \in \mathcal{D}_{\omega_i}$ will inherit the access level label of the respective data item, i.e. $\lambda(id_{d_i}) = \lambda(d_i)$.

Alg. 1 $(K_O, K_S, PP) \stackrel{\$}{\leftarrow} \text{KeyGen}(1^\kappa, \mathcal{U}, \mathcal{P})$

```

1: for  $i \in |\mathbb{L}|$  do
2:    $k_{a_i,1}, k_{a_i,2}, k_{a_i,3} \stackrel{\$}{\leftarrow} \{0, 1\}^\kappa$ 
3:    $(PP_{\text{BE}}, k_{\text{BE}}) \stackrel{\$}{\leftarrow} \text{BE.KeyGen}(1^\kappa, |\mathcal{U}|)$ 
4:    $st_O \stackrel{\$}{\leftarrow} \{0, 1\}^\kappa$ 
5:    $S \stackrel{\$}{\leftarrow} \mathcal{U}$ 
6:    $\mathcal{G} \leftarrow \{S\}$ 
7:    $st_S \stackrel{\$}{\leftarrow} \text{BE.Enc}(st_O, \mathcal{G}, k_{\text{BE}})$ 
8: return
9:  $K_O \leftarrow (\{k_{a_i,1}\}_{i \in |\mathbb{L}|}, \{k_{a_i,2}\}_{i \in |\mathbb{L}|}, \{k_{a_i,3}\}_{i \in |\mathbb{L}|}, k_{\text{BE}}, st_O)$ 
10:  $PP \leftarrow (\mathcal{P}, \mathcal{G}, st_S, PP_{\text{BE}})$ 
11:  $K_S \stackrel{\$}{\leftarrow} \text{BE.Add}(k_{\text{BE}}, S)$ 

```

Alg. 2 $\mathcal{I}_{\mathcal{D}} \stackrel{\$}{\leftarrow} \text{BuildIndex}(\mathcal{D}_{aug}, K_O, PP)$

```

1:  $free \leftarrow \{\text{addr}(\mathbb{N}_i)\}_{i \in |\mathbb{A}|}$ 
2: for  $1 \leq i \leq |\mathcal{W}|$  do
3:    $(\mathbb{L}_{\omega_i}, \mathbb{X}_{\omega_i}, \mathbb{N}_{\omega_i}) \leftarrow \gamma(\mathcal{D}_{\omega_i})$ 
4:    $free \leftarrow free \setminus \mathbb{N}_{\omega_i}$ 
5:   for  $1 \leq j \leq |\mathbb{N}_{\omega_i}| - 1$  do
6:      $r_j \stackrel{\$}{\leftarrow} \{0, 1\}^\kappa$ 
7:      $\mathbb{A}[\mathbb{N}_{\omega_i}[j]] \leftarrow \left( (\mathbb{L}_{\omega_i}[j], \mathbb{N}_{\omega_i}[j+1], P_{k_{\lambda(\mathbb{L}_{\omega_i}[j+1]),3}}(\omega_i)}) \oplus \right.$ 
8:        $\left. H(P_{k_{\lambda(\mathbb{L}_{\omega_i}[j]),3}}(\omega_i), r_j), r_j \right)$ 
9:      $\mathbb{A}[\mathbb{N}_{\omega_i}[|\mathbb{N}_{\omega_i}|]] \leftarrow \left( (\mathbb{L}_{\omega_i}[|\mathbb{N}_{\omega_i}|], 0, 0) \oplus \right.$ 
10:        $\left. H(P_{k_{\lambda(\mathbb{L}_{\omega_i}[|\mathbb{N}_{\omega_i}|]),3}}(\omega_i), r_{|\mathbb{N}_{\omega_i}|}), r_{|\mathbb{N}_{\omega_i}|} \right)$ 
11:   for  $1 \leq \ell \leq |\mathbb{L}|$  do
12:     if  $\mathbb{X}_{\omega_i}[a_\ell] \neq \perp$  then
13:        $\mathbb{T}[F_{k_{a_\ell,1}}(\omega_i)] \leftarrow (\mathbb{N}_{\omega_i}[\mathbb{X}_{\omega_i}[a_\ell]] \oplus G_{k_{a_\ell,2}}(\omega_i))$ 
14:     else
15:        $\mathbb{T}[k_{a_\ell,1}(\omega_i)] \leftarrow \perp$ 
16:   return  $\mathcal{I}_{\mathcal{D}} \leftarrow (\mathbb{A}, \mathbb{T})$ 

```

Alg. 3 $(K_u, PP) \stackrel{\$}{\leftarrow} \text{AddUser}(u, \lambda(u), K_O, PP)$

```

1:  $\mathcal{G} \leftarrow \mathcal{G} \cup \{u\}$ 
2:  $k_u \stackrel{\$}{\leftarrow} \text{BE.Add}(k_{\text{BE}}, u)$ 
3:  $st_S \stackrel{\$}{\leftarrow} \text{BE.Enc}(st_O, \mathcal{G}, k_{\text{BE}})$ 
4: return
5:  $PP \leftarrow (\mathcal{P}, \mathcal{G}, st_S, PP_{\text{BE}})$ 
6:  $K_u \leftarrow (k_u, k_{\lambda(u),1}, k_{\lambda(u),2}, k_{\lambda(u),3})$ 

```

Alg. 4 $(K_O, PP) \stackrel{\$}{\leftarrow} \text{RevokeUser}(u, K_O, PP)$

```

1:  $st_O \stackrel{\$}{\leftarrow} \{0, 1\}^\kappa$ 
2:  $\mathcal{G} \leftarrow \mathcal{G} \setminus \{u\}$ 
3:  $st_S \stackrel{\$}{\leftarrow} \text{BE.Enc}(st_O, \mathcal{G}, k_{\text{BE}})$ 
4: return  $K_O \leftarrow (\{k_{a_i,1}\}_{i \in [|\mathbb{L}|]}, \{k_{a_i,2}\}_{i \in [|\mathbb{L}|]}, \{k_{a_i,3}\}_{i \in [|\mathbb{L}|]}, k_{\text{BE}}, st_O)$ 

```

Alg. 5 $(T_{\omega,\lambda(u)} \leftarrow \text{Query}(\omega, K_u))$

```

1:  $st'_O \leftarrow \text{BE.Dec}(k_u, st_S)$ 
2: if  $st'_O = \perp$  then
3:   return  $\perp$ 
4:  $t_{\omega,\lambda(u)} \leftarrow (F_{k_{\lambda(u),1}}(\omega), G_{k_{\lambda(u),2}}(\omega), P_{k_{\lambda(u),3}}(\omega))$ 
5: return
6:  $T_{\omega,\lambda(u)} \leftarrow \phi_{st'_O}(t_{\omega,\lambda(u)})$ 

```

Alg. 6 $(\mathcal{R}_{\omega,\lambda(u)} \leftarrow \text{Search}(T_{\omega,\lambda(u)}, \mathcal{I}_D, K_S))$

```

1:  $st'_O \leftarrow \text{BE.Dec}(K_S, st_S)$ 
2: Parse  $\phi_{st'_O}^{-1}(T_{\omega,\lambda(u)})$  as  $(\tau_1, \tau_2, \tau_3)$ 
3:  $\mathcal{R}_{\omega,\lambda(u)} \leftarrow \emptyset$ 
4: if  $\mathbb{T}[\tau_1] = \perp$  then
5:   return  $\perp$ 
6:  $v_2 \leftarrow 1$ 
7: while  $v_2 \neq 0$  do
8:   Parse  $\mathbb{T}[\tau_1] \oplus \tau_2$  as  $y$ 
9:   Parse  $\mathbb{A}[y]$  as  $(z_1, z_2)$ 
10:  Parse  $z_1 \oplus H(\tau_3, z_2)$  as  $(v_1, v_2)$ 
11:   $\mathcal{R}_{\omega,\lambda(u)} \leftarrow \mathcal{R}_{\omega,\lambda(u)} \cup \{v_1\}$ 
12: return  $\mathcal{R}_{\omega,\lambda(u)}$ 

```

The **KeyGen** algorithm initialises the system and generates the keys K_O, k_S , along with the public parameters, PP. The key K_O includes the secret key for the BE scheme and the sets of $|\mathbb{L}|$ keys for each pseudo-random function: F, G and P and the key for the pseudo-random permutation ϕ (referred to as the data owner's state, st_O). The server is enrolled as a user and its secret key is also generated (although it does not receive the necessary keys to form search queries). PP includes the information flow policy \mathcal{P} , the authorized user group \mathcal{G} , the server state st_S (which is an encryption of the owner state generated using BE) and the public parameters for BE, PP_{BE} .

The **BuildIndex** algorithm initializes a set $free$ which consists of all nodes in the array \mathbb{A} . **BuildIndex** considers each keyword contained in the dataset in turn. For each keyword ω_i , the function γ generates $L_{\omega_i}, X_{\omega_i}$ and N_{ω_i} . The free list is then updated according to which nodes have been chosen by γ . The nodes in the array that form the linked lists consist of the identifier from L_{ω_i} of a data item containing ω_i , the address in the array of the next node in the linked list, the key used to decrypt the following node in the linked list and a random bit string $r_i \in \{0, 1\}^k$. The identifier, address of the next node and the key used to decrypt the following node in the linked list are XORed with the output of a PRF H in order to encrypt this information. For the first node in the linked list the input of H is the decryption key for the current node (which corresponds to an access level and keyword and forms part of the search query) along with r_i , hence the information stored in the node can only be decrypted by the server if the server has a search query generated by a user who is authorized to view the data item whose identifier is stored at that node. The decryption key for all subsequent nodes is contained in the previous node of the linked list. **BuildIndex** then proceeds to create the lookup table \mathbb{T} . Unlike prior schemes [9], each user may have a different access level and thus the starting points for search results within the linked lists may vary; a search query made

by a user with a higher access level should traverse more of the list than that of a user with lower access rights (the user is authorised to search more data items). Table \mathbb{T} has an entry for each access level/keyword pair containing the address of a node in \mathbb{A} , which is the node in the linked list L_{ω_i} from which the user with a specified access level is authorised to decrypt. If an access level is not authorised to view any part of the linked list then the value in \mathbb{T} is set to \perp . Finally the index $\mathcal{I}_D = (\mathbb{A}, \mathbb{T})$ is returned.

The **AddUser** algorithm grants a user u the ability to search the index at a specific access level. The user is added to the set \mathcal{G} of authorized users and a BE key, k_u , is derived for the new user. The new user is given k_u and the secret keys associated with their access level $k_{\lambda(u),1}, k_{\lambda(u),2}$ and $k_{\lambda(u),3}$ and PP is updated.

The **RevokeUser** algorithm revokes a user's search privileges. A new value for st_O is selected and the user is removed from \mathcal{G} . This value is encrypted using BE to form the new server state st_S . The updated versions of K_O and PP are output.

The **Query** algorithm generates a search query for a user u to search for a keyword w . The user first attempts to decrypt the current server state st_S using their secret key k_u ; we denote the output of the decryption by st'_O . Note that if u is not authorised then decryption will return \perp , if this is the case **Query** outputs \perp . The query itself comprises three parts. The first is the output of the PRF F applied to the keyword ω , keyed with the secret key for F associated with the user's access level $k_{\lambda(u),1}$. This part of the query is used to locate the relevant entry in \mathbb{T} . The second part is the output of the PRF G applied to the keyword ω and is used to mask the entry in \mathbb{T} in order to locate the user's relevant starting position in the linked list corresponding to ω in \mathbb{A} . The third part is the output of the PRF P applied to the keyword ω , which is used to decrypt the first relevant node in \mathbb{A} according to the user's access level. The PRP ϕ is applied to the search query, using st'_O as the key.

The **Search** algorithm finds data item identifiers associated with the searched keyword from the subset of data item identifiers the user is authorized to search. The server decrypts st_S and applies the inverse of the PRP ϕ to the query it received; it parses the result as (τ_1, τ_2, τ_3) . The server then looks up entry $\mathbb{T}[\tau_1]$ and if that entry is not equal to \perp , the server XORs the value with τ_2 and parses the resulting value as y . The server looks up the node at $\mathbb{A}[y]$, parses the entry as (z_1, z_2) , and decrypts it by XORing z_1 with the output of H (which takes as input τ_3 along with z_2).

The server is able to sequentially decrypt the rest of the list stored in \mathbb{A} until they reach a node where the address stored in that node for the next item in the linked list is 0.

3.4 Security

In MLSSE search queries for the same keyword that are produced by users with different access levels are indistinguishable from one another. That is, a search query for a keyword ω from a user u_i with access level $\lambda(u_i)$ is indistinguishable from a search query for ω from a user u_j with access level $\lambda(u_j)$ for $\lambda(u_i) \neq \lambda(u_j)$. This means that from the queries alone an adversary is unable to deduce how many times a certain keyword has been searched for overall, it can only deduce how many times the same keyword has been searched for within each access level. This information leakage is less than that of standard single or multi user SSE

schemes such as [6, 7, 9, 10, 12, 13].

In terms of access pattern we also reduce the amount of information leakage compared with standard single user or multi-user SSE schemes. In particular we do not reveal whether a data item contains the keyword ω_i associated with a search query unless the access level of that data item is less than or equal to that of the user u_i that generated the search query, meaning that an adversary cannot see a full set of search results.

However when a search query is paired with the search results it generates (the access pattern, Definition 5) then an adversary may be able to correlate which search queries are for the same keyword by looking at the intersections of the search results. For example if one set of search results is a subset of another set of search results then this may imply that the two search queries used to generate these results are for the same keyword. An adversary may eventually be able to build up a complete set of search results for a particular keyword, which is equivalent to the leakage produced by a search query in a single user SSE scheme. The server does not know, however, how many access levels there are altogether so a server would need to receive all possible search queries before it can ascertain whether or not a set of search results for a particular keyword is complete or not.

The hierarchical relationships between the data item identifiers i.e. which identifiers represent data items at higher access level than others could also be leaked in the same way. If an adversary has ascertained that two sets of search results $\mathcal{R}_{\omega, a_i} \subset \mathcal{R}_{\omega, a_j}$ represent searches for the same keyword ω , then an adversary will be able to conclude that identifiers in the set $\mathcal{R}_{\omega, a_j} \setminus \mathcal{R}_{\omega, a_i}$ are at a higher access level than those in $\mathcal{R}_{\omega, a_i}$. We note that unless the search results are padded in some way this leakage is inevitable. Padding search results is not standard in SSE schemes as it requires post-processing of the search results by the user hence we do not pad the search results in our system model in order to maintain an efficient scheme.

From this we can see that initially our scheme leaks less information about the search pattern and access pattern than a single user SSE scheme, however over time as more queries are generated the information leakage tends to that of a single user SSE scheme. The information leakage relating to a keyword ω i.e. the access patterns for search queries corresponding to ω only reaches that of a single user SSE scheme once a search query has been generated at each possible access level, our leakage remains lower up until this point.

As a search query for a keyword and access level pair is created deterministically we can think of the search query as a *codeword* for the combination of that keyword and access level. The index usually reveals these codewords as a search is carried out by matching search queries to relevant codewords in the index. A codeword for keyword ω at access level a is denoted $id(\omega, a)$.

We give the specific leakage functions to precisely capture the leakage in MLSSE, where Ω is a set of queries from users in the system that have been evaluated on the encrypted index by the server:

1. $\mathcal{L}_{Setup}(\mathcal{I}_D) = (|\mathbb{A}|, |\mathbb{T}|, [id(\omega, a)]_{\omega \in \Delta, i \in [|\mathbb{L}|]})$
2. $\mathcal{L}_{Query}(\mathcal{I}_D, \Omega) = (AP(\mathcal{I}_D, \Omega), SP(\mathcal{I}_D, \Omega), [id(\omega, a)]_{\forall T_{\omega, a} \in \Omega, \Omega})$

THEOREM 1. *Given an IND-CPA secure broadcast encryption scheme BE, a pseudo-random permutation ϕ , and pseudorandom functions F, G, P, H . Let MLSSE be the searchable*

symmetric encryption scheme with multi-level access defined in Algorithms 1-6. Then MLSSE is $(\mathcal{L}_{Setup}, \mathcal{L}_{Query})$ -secure in the sense of multi-level access and revocation.

We provide the intuitions of our security proofs here and refer the reader to the full online version of the paper for the full security proofs [1].

Multi-level access: To show multi-level access we reduce the security to that of the IND-CPA security of a symmetric encryption scheme which encrypts plaintexts by XORing them with the output of a PRF. We assume the possibility of an adversary \mathcal{A} that is able to break the multi-level security of our scheme then we construct a second adversary \mathcal{A}' that is able to use \mathcal{A} as a subroutine in order to break the IND-CPA security of the symmetric encryption scheme with non-negligible probability.

Revocation: In this proof we show that if we assume an adversary \mathcal{A} with non-negligible advantage δ in Game 2 then \mathcal{A} can be used as a subroutine by an adversary \mathcal{A}_{BE} to break the security of an IND-CPA secure broadcast encryption scheme BE.

3.5 Achieving dynamicity

We can extend MLSSE to support multi-level access on a dynamic data set by adding two new data structures to the index: a deletion table (\mathbb{T}_d) and a deletion array (\mathbb{A}_d). There are also four additional algorithms: **AddToken**, **Add**, **DeleteToken**, **Delete**. Array \mathbb{A}_d stores a list of nodes for each data item which point to nodes in \mathbb{A} that would need to be removed if the corresponding data item was deleted. This means that every node in \mathbb{A} will have a corresponding node in \mathbb{A}_d , which is called its *dual* node. \mathbb{T}_d is a table with an entry for each data item which points to the start of the corresponding linked list in \mathbb{A}_d , given a valid *delete token* for that data item. In addition to these two new structures the index consists of a search array \mathbb{A}_s and a search table \mathbb{T}_s (as in the original construction) and a *free list* that keeps track of all the unused space in \mathbb{A}_s .

In the dynamic scheme searching for a keyword is done similarly to the static construction in Section 3.3 and follows the concept of linked lists presented by [9].

To add a data item to the index, changes need to be made to $\mathbb{T}_d, \mathbb{A}_s$ and \mathbb{A}_d . The data owner creates an *add token* using **AddToken** and sends this to the server. The server then determines the free space available in \mathbb{A}_s using the free list and adds the relevant information to the free nodes and updates the free list. When adding a new data item the relevant nodes cannot be added to the end of each linked list; instead we have to insert in the appropriate place in the linked list according to the access level of the new data item. Information in the add token will allow the server to locate the correct point at which to insert the nodes in each linked list, so instead of the entry in \mathbb{T}_s just pointing to the end node of each linked list this is altered so that it points to the correct node in the linked list according to the access level of the new data item. The respective predecessor of each new node is modified to point to the new node instead of its previous ancestor.

In order to remove a data item, a deletion token is created which allows the server to locate and delete the correct entries in \mathbb{T}_d . This, in turn, allows the server to locate and delete the correct entries in \mathbb{A}_s . Some nodes will need to be updated in \mathbb{A}_s (as some of the linked lists will have nodes which point to nodes that have been deleted) and this is

done using homomorphic encryption.

3.6 Efficiency

In this section we discuss the efficiency of our multi-user, multi-level construction compared with the single-user construction of [13]. As our scheme is static and the scheme of [13] is dynamic, we ignore the structures and algorithms in [13] that apply to the dynamicity, such as the deletion table, the deletion array and algorithms `AddToken`, `Add`, `DeleteToken`, `Delete`.

The index is composed of a look-up table and a search array. No changes are made to the search array that effect the time needed to generate it or the search time, but the look-up table needs to be augmented by a factor of $|\mathbb{L}|$; this will require more space on the server but does not effect the search time. The size of our index is $\mathcal{O}(\Delta \cdot |\mathbb{L}| + n)$ whereas the size of the index in the single user scheme is $\mathcal{O}(\Delta + n)$.

There search time of our scheme is $\mathcal{O}(|\mathcal{D}_{\omega,a}|)$ where $\mathcal{D}_{\omega,a}$ is the set of data item identifiers satisfying the search query $T_{\omega,a}$. This is equivalent to the search time of [13], however in our scheme the size of $\mathcal{D}_{\omega,a}$ is likely to be smaller, depending on the access level of the user who generated the search query.

The amount of computation required to generate the search queries as well as the size of the search queries is the same in both schemes, they are both constructed by evaluating three PRFs.

We note that in terms of efficiency our construction is very similar to that of [13]. This is also true for the dynamic version of our construction.

4. CONCLUSION

We have defined a new system, security models and a construction for symmetric solutions to searching on encrypted data in the multi-level setting. Users may search for keywords within a set of encrypted data items, restricting the search to data items they are authorised to view only. Future work will focus on increasing the range of query types beyond that of single keyword equality search and to expand the access control policies to arbitrary information flow policies.

5. REFERENCES

- [1] J. Alderman, K. M. Martin, and S. L. Renwick. Multi-level access in searchable symmetric encryption. IACR Cryptology ePrint Archive, Report 2017/211, 2017.
- [2] E. Bell and L. La Padula. Secure computer system: Unified exposition and multics interpretation. Technical report, Mitre Corporation, 1976.
- [3] J. Benaloh, M. Chase, E. Horvitz, and K. E. Lauter. Patient controlled encryption: ensuring privacy of electronic medical records. In *Proceedings of the first ACM Cloud Computing Security Workshop, CCSW 2009*, pages 103–114. ACM, 2009.
- [4] D. Boneh, G. Di Crescenzo, R. Ostrovsky, and G. Persiano. Public key encryption with keyword search. In *Advances in Cryptology - EUROCRYPT 2004, International Conference on the Theory and Applications of Cryptographic Techniques*, volume 3027 of *Lecture Notes in Computer Science*, pages 506–522. Springer, 2004.
- [5] J. W. Byun, H. S. Rhee, H. Park, and D. H. Lee. Off-line keyword guessing attacks on recent keyword search schemes over encrypted data. In *Secure Data Management, Third VLDB Workshop, SDM 2006*, volume 4165 of *Lecture Notes in Computer Science*, pages 75–83. Springer, 2006.
- [6] Y. Chang and M. Mitzenmacher. Privacy preserving keyword searches on remote encrypted data. In *Applied Cryptography and Network Security, Third International Conference, ACNS 2005*, volume 3531 of *Lecture Notes in Computer Science*, pages 442–455. Springer, 2005.
- [7] M. Chase and S. Kamara. Structured encryption and controlled disclosure. In *Advances in Cryptology - ASIACRYPT 2010 - 16th International Conference on the Theory and Application of Cryptology and Information Security*, volume 6477 of *Lecture Notes in Computer Science*, pages 577–594. Springer, 2010.
- [8] J. Crampton. Cryptographic enforcement of role-based access control. In *Formal Aspects in Security and Trust*, volume 6561 of *Lecture Notes in Computer Science*, pages 191–205. Springer, 2010.
- [9] R. Curtmola, J. A. 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, CCS 2006*, pages 79–88. ACM, 2006.
- [10] E.-J. Goh. Secure indexes. IACR Cryptology ePrint Archive, Report 2003/216, 2003.
- [11] A. Kaci, T. Bouabana-Tebibel, and Z. Challal. Access control aware search on the cloud computing. In *2014 International Conference on Advances in Computing, Communications and Informatics, ICACCI 2014*, pages 1258–1264. IEEE, 2014.
- [12] S. Kamara and C. Papamonthou. Parallel and dynamic searchable symmetric encryption. In *Financial Cryptography and Data Security - 17th International Conference, FC 2013*, volume 7859 of *Lecture Notes in Computer Science*, pages 258–274. Springer, 2013.
- [13] S. Kamara, C. Papamonthou, and T. Roeder. Dynamic searchable symmetric encryption. In *The ACM Conference on Computer and Communications Security, CCS'12*, pages 965–976. ACM, 2012.
- [14] Z. A. Kissel and J. Wang. Verifiable symmetric searchable encryption for multiple groups of users. In *Proceedings of the 2013 International Conference on Security and Management*, pages 179–185. CSREA Press, 2013.
- [15] M. Li, S. Yu, N. Cao, and W. Lou. Authorized private keyword search over encrypted data in cloud computing. In *2011 International Conference on Distributed Computing Systems, ICDCS*, pages 383–392. IEEE Computer Society, 2011.
- [16] C. Office. Government security classifications. Technical report, 2013.
- [17] M. Onen, R. Molva, and C. Van Rompay. Multi-user searchable encryption in the cloud. In *Information Security - 18th International Conference, ISC 2015*, volume 9290 of *Lecture Notes in Computer Science*, pages 299–316. Springer, 2015.

- [18] D. X. Song, D. Wagner, and A. Perrig. Practical techniques for searches on encrypted data. In *2000 IEEE Symposium on Security and Privacy*, pages 44–55. IEEE, 2000.
- [19] W. Sun, S. Yu, and W. Lou. Protecting your right: Attribute-based keyword search with fine-grained owner-enforced search authorization in the cloud. In *2014 IEEE Conference on Computer Communications, INFOCOM 2014*, pages 226–234. IEEE, 2014.
- [20] W. Sun, S. Yu, W. Lou, T. Hou, and H. Li. Protecting your right: Verifiable attribute-based keyword search with fine-grained owner-enforced search authorization in the cloud. *IEEE Transactions on Parallel Distributed Systems*, 27(4):1187–1198, 2016.
- [21] Y. Yang. Attribute-based data retrieval with semantic keyword search for e-health cloud. *Journal of Cloud Computing: Advances, Systems and Applications*, 4, 2015.