

Storage Efficient Substring Searchable Symmetric Encryption

Iraklis Leontiadis, Ming Li

University of Arizona, USA
{leontiad,lim}@email.arizona.com

Abstract. We address the problem of substring searchable encryption. A single user produces a big stream of data and later on wants to learn the positions in the string that some patterns occur. Although current techniques exploit auxiliary data structures to achieve efficient substring search on the server side, the cost at the user side may be prohibitive. We revisit the work of substring searchable encryption in order to reduce the storage cost of auxiliary data structures. Our solution entails suffix array which allows optimal storage cost $O(n)$ with small hidden factor at the size of the string n . On top of that we build an encrypted index that allows the server to answer substring queries without learning neither the query nor the result. We identify the leakages of the scheme following the work of Curtmola *et al.* [12] and we analyze the security of the protocol in the real ideal framework. Moreover, we demonstrate the practicality of the protocol by searching a one million characters data stream in less than one second within the GPU computing paradigm. The total speedup approximates a factor of 4x, compared with naive CPU implementation.

1 Introduction

Nowadays there is a flourish of protocols delegated to run by an untrusted coalition of servers, systems, services, called hereafter the cloud. Due to the untrusted nature of the cloud, users seek to protect the privacy and security of their data with cryptographic primitives. The cloud on the other hand offers an economy of scale with the impressive resources it acquires, ranging from software to hardware. Uploading encrypted data however, renders operation on it infeasible. Downloading, decrypting and running the operation on plaintext data cancels the advantages that the cloud offers for large storage and computational efficiency. Usually users need to perform a search on their data. Tailored protocols for secure searchable encryption have been proposed in the literature, whereby single or multiple users upload encrypted documents, with some auxiliary data structure called an index, allowing the cloud to correctly return documents containing a single, multiple or a boolean function of keywords, without compromising index, query, and documents privacy. Apart from their theoretical consideration in the literature, quite a few companies adopt this model to offer searchable encryption schemes over encrypted data [3, 10, 11, 25, 32, 38, 42].

While keyword based search protocols are quite common in a large range of applications, they cannot efficiently address all the possible queries a user submits to the cloud. Substring based queries have come to the forefront due to the ubiquitousness of devices and the progress in storage technology. Devices produce a big stream of data, which needs to be queried later on with *substring based queries*. Namely a substring query for a stream of data, consists of a substring of the stream and the result is the position of the substring in the big stream, or/and the number of occurrences of multiple substrings.

The variety of application for *substring based queries* spans in financial applications, health-care analysis and surveillance of malicious behavior. In surveillance based applications an authority through espionage, logs chat rooms of suspicious conversations for criminal activities. The goal of the authority is to identify the involved chat room participants. Suspicious parties at some point post identifiable information such as secret url pages. Authority keeps track of the pages and coerces the cloud service through jurisdictional power to reveal if this specific url (substring) obtained by a chat room is part of their logs in order to identify the suspicious participants. The source of information could also be a unique identifiable transaction on a cryptocurrency based ledger. The authority has identified a criminal based transaction and needs to obtain further information for the transaction (its positions on the ledger) in order to proceed on further investigations. In a health-care application, data enclaves which hold giant stream of medical information such as DNA sequencing are asked to answer substring queries by medical labs. The possible position of a substring in the whole DNA sequence of a single person gives information about predisposition to deceases.

Protecting the privacy of the data stream and the substring query, while allowing an untrusted cloud to correctly answer substring matching pattern efficiently and securely is not trivial. Following the searchable encryption approach, separating the data itself from the index, results in a prohibitive storage index cost $O(n^2)$, where n is the size of the stream. The index would consist of all the possible substrings n of a stream of data of size n and the encrypted data would be the positions of the substring. Recently Chase *et al.* [9] proposed a solution that asymptotically achieves $O(n)$ storage costs by exploiting the auxiliary data structure of the suffix tree. However the asymptotic costs of $O(n)$ hide a constant factor that can be roughly up to $20n$ [1, 5, 22, 30] for the construction of the suffix tree due to the complexity of the tree and the extra pointers to traverse a tree. Moreover the suffix tree based approach leaks unnecessary that eventually can reveal all the encrypted positions of the substrings. Following a different approach other than auxiliary index based methods, the authors in [16] achieve to hide the extra leakage at the cost of fixed length substring patterns. The neat of their solution lies on the design of subset sum problems tailored to the positions of specific substrings, so as to the cloud can solve it partially. However this comes at the cost of small constant substring query length during the execution of the protocol.

In this paper we design and analyze a storage efficient Substring Searchable Symmetric Encryption (S^3E) protocol with minimal leakage and variable size of substrings. We follow a different approach of existing techniques that allows us to achieve the efficiency, functional and security goals we want. In our technique we exploit a self-indexed data structure, which allows the cloud to search for substring queries without the need for the original stream of data. Its form resembles the suffix arrays with some additional extra steps, thus we are taking its computational cost for “free”, after building the suffix array. From the index the cloud can answer substring queries efficiently. The main contributions in the paper are summarized as follows:

- *Storage efficient Substring Searchable Symmetric Encryption (S^3E):* Thanks to the employment of the suffix array, which achieves a small hidden factor (≈ 4) in the $O(n)$ asymptotic complexity, compared to the bigger (≈ 20) hidden factor of the suffix tree, our design presents a storage efficient substring searchable symmetric encryption protocol.
- *Minimum leakage:* The simplicity of the suffix array construction hides unnecessary information that is leaked when a suffix tree is used, such as the visited leaves of the tree between various queries.
- *Variable substring query length:* Our solution allows a dynamic issue of substring queries of variable size without the need of defining a fixed query size beforehand.
- *On/Off-line Difference attack:* Throughout our novel design principle based on suffix arrays and a self-indexed data structure which encompasses both the index and the data, we identified a novel attack called hereafter on/off-difference attack. An adversary who observes the encrypted index and the query response is able to fully decrypt the encrypted positions with the differences it learns while traversing the index. Finally, we mitigate the on/off-line difference attack by encrypting the addresses of the suffix array instead of the elements thereof and running a PIR protocol in the final round of the protocol.
- *Prototype implementation:* We implement our protocol in order to demonstrate its practicality. We use off the shelf laptops for the implementation of the index, which is done at the client side. Moreover, we take advantage of the multiple available cores in a GPU environment, which allow parallelization of costly computations at the cloud side during a query response, with increased computation to memory ratio.

Outline. In section 2 we introduce the problem this paper addresses. Afterwards, in section 3 we review similar cryptographic protocols for substring searchable symmetric encryption. We continue in section 4 with the core idea of our solution. The full protocol description is presented in section 5. We then investigate the security and the costs of the proposed scheme for storage efficient substring searchable symmetric encryption. In section 6 we present our prototype implementation results. Finally, we conclude in section 7.

2 Problem Statement

In this section we formalize first the problem of string matching. We first start with the functional requirements of substring matching and afterwards we present the security requirements of the protocol.

2.1 Functional Requirements

Herewith pattern matching, string matching and substring matching are used interchangeably in this paper. We assume that a string S is modeled as an one dimension array $S[1..n]$. A substring is another array $T[1..m]$. The elements of each array are drawn from some finite alphanumerical alphabet Σ of size $\sigma = |\Sigma|$. We say that a substring T occurs in S if there exists $s : 1 \leq s \leq n - m$ and $T[s + 1..s + m] = S[1..m]$, meaning that $S[s + j] = T[j], 1 \leq j \leq m$ (cf figure 1). If such position in S does not exist the algorithm returns \perp . A string x is prefix of a string $y : x \prec y$ if $y = xw$ for some $w \in \Sigma$. Similarly a string x is a suffix of string $y : x \succ y$ if $y = wx$ for some string $w \in \Sigma$.

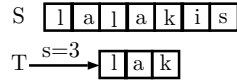


Fig. 1: Pattern matching

Naive algorithms for pattern matching achieve $O(n)$ on search time and 0 cost on preprocessing. The algorithm simply scans all the positions $i, 1 \leq i \leq n - m$ of the string S as long as it finds m consecutive matches at a position $j, 1 \leq j \leq n - m + 1$. Trading preprocessing efficiency for better search costs, Robin Karp algorithm [27] achieves $O(n - m + 1)$ search time and $\Theta(m)$ preprocessing amortized cost. In a similar trajectory Knuth-Morris-Pratt [28] has $\Theta(n)$ search complexity and $\Theta(m)$ preprocessing time. Boyer-Moore pattern matching technique [4] increases the preprocessing cost at $\Theta(m + \sigma)$ in order to have worst case search complexity $O(n)$. Following a different trajectory substring matching techniques achieve $O(m)$ search time by leveraging a more sophisticated preprocessing step, in which the suffixes of all substring are computed along with their positions in the string S , be it suffix tree [31, 41, 43] or suffix array [30]. Suffix tree though has a more expensive space efficiency due to the extra information the suffix tree has to keep [1, 5, 22, 30]. This cost is translated to a constant factor that approximates ≈ 20 , which is hidden in the $O(n)$ asymptotic storage cost of the suffix tree construction. As a first step to relax this storage extra hidden cost we choose to build upon the suffix array string matching approach which has a much simpler storage cost which approximates $4n$ [1].

We redraw upon the *queryable encryption* syntactical definition of [9], since we believe it follows a deceptive abstraction. Namely, the functional definition claims to capture a generic framework for searchable encryption, in the sense that a query \mathcal{F} can be any function keyword query, or substring query. However, an encrypted searchable encryption scheme is a more generic protocol, since it can be used to solve the substring searchable encryption problem with the encrypted inverted index technique as shown in the introduction. As such searchable and substring encryption schemes cannot be addressed by the same definitional framework, but most importantly. Furthermore, the nature of the problem and the solution for substring queries drastically varies from keyword searchable encryption, since the index contains the data and there are not two separate objects, meaning that the index for substring queries is *self-indexed*, since from the index you can recover the underlying data structure. In contrast in encrypted searchable encryption, there is a clear distinction between the index, and the data structure that holds the data (files with keywords). For these reasons we rewrite the functional definitional framework for substring searchable encryption in order to capture with more precision its semantics.

Definition 1. A *Substring Searchable Symmetric Encryption scheme* (S^3E) is a collection of four polynomial time algorithms ($\text{KeyGen}, \text{PreProcess}, \text{SrchToken}, \text{Search}$) defined as follows:

- $k \leftarrow \text{KeyGen}(1^\lambda)$: It is a probabilistic algorithm that takes as input the security parameter in the unary form 1^λ and outputs the secret substring search key k .
- $\text{SES} \leftarrow \text{PreProcess}(k, S)$: This algorithm takes as inputs the stream S and the secret key k and outputs the substring encrypted data structure SES .
- $\text{tk}_{T,S} \leftarrow \text{SrchToken}(k, T[1..m])$: It is a probabilistic algorithm that takes as input the secret substring search key k , a string $T[1..m]$ and outputs a trapdoor to search for the string T on data stream S , through SES .
- $(s, \perp) \leftarrow \text{Search}(\text{tk}_{T,S}, \text{SES})$: It is a deterministic algorithm which takes as input a trapdoor $\text{tk}_{T,S}$ and an substring encrypted structure SES and outputs the positions s in S that substring T occurs, or \perp otherwise.

A S^3E is correct if $\forall \lambda \in \mathbb{N}, \forall S \in \Sigma, \forall k \leftarrow \text{KeyGen}(1^\lambda), \forall \text{SES} \leftarrow \text{PreProcess}(k, S), \forall \text{tk}_{T,S} \leftarrow \text{SrcToken}(k, T[1\dots m]), \text{Search}(\text{tk}_{T,S}, \text{SES})$ always returns the correct positions s in the string S or \perp otherwise.

2.2 Security Model

Intuitively the security guarantee we ask for is **1**) given a probabilistic polynomial time adversary \mathcal{A} with access to a substring encrypted structure SES , \mathcal{A} cannot gain more partial information about the underlying stream of data S and **2**) given a set of trapdoor tokens for an adaptive generated set of queries $\mathbf{q} = (q_1, q_2, q_3, \dots, q_o)$ associated with set of tokens $\mathbf{t} = (\text{tk}_1, \text{tk}_2, \text{tk}_3, \dots, \text{tk}_o)$ \mathcal{A} cannot learn anything for \mathbf{q} and \mathbf{t} . Following the symmetric searchable paradigm we know it is impossible to achieve those two security guarantees without leaking some extra information as the observed in [7, 8, 12].

We express the security guarantees of the protocol in terms of the consolidated simutability [29]. First a leakage function \mathcal{L} is defined, which expresses the leakage of a S^3E scheme to an adversary \mathcal{A} , though the transcripts of the protocol. The simulation framework assumes two games. The $\mathbf{Real}_{\mathcal{A}(\lambda)}^{S^3E}$ game, in which adversaries can corrupt the parties they want and the $\mathbf{Ideal}_{\mathcal{A}, \mathcal{S}(\lambda)}^{S^3E}$ one in which there is only benign behavior of each party. The security analysis narrows down to the design of a simulator \mathcal{S} , who tries to simulate the malicious behavior in the $\mathbf{Ideal}_{\mathcal{A}, \mathcal{S}(\lambda)}^{S^3E}$ game, only through access to the leakage function \mathcal{L} . We say that a protocol is secure if \mathcal{S} simulates indistinguishable views of the adversary \mathcal{A} in the $\mathbf{Ideal}_{\mathcal{A}, \mathcal{S}(\lambda)}^{S^3E}$ game.

Definition 2. A leakage function \mathcal{L} for a S^3E scheme comprises the following three leakage functions:

- PreProcess Leakage: \mathcal{L}_1 includes the padded size of the data stream $n' \geq |S|$.
- SrcToken Leakage: The SrcToken Leakage \mathcal{L}_2 reveals the length of the token $|\text{tk}|$ and how many common characters reside in it.
- Search Leakage: \mathcal{L}_3 leaks the how many times a substring token tk exists in string S .

The adversary \mathcal{A} plays the role of a semi-honest cloud and during the two games we assume a challenger \mathcal{C} who interacts with \mathcal{A} . We describe the two games in algorithmic details in what is follows: $\mathbf{Real}_{\mathcal{A}(\lambda)}^{S^3E}$ game:

- \mathcal{C} runs $\text{KeyGen}(1^\lambda)$ to obtain k .
- \mathcal{A} chooses a string $S \in \Sigma$, sends it to \mathcal{C} and \mathcal{C} replies with $\text{SES} \leftarrow \text{PreProcess}(k, S)$ to \mathcal{A} .
- \mathcal{A} issues a polynomial number of adaptively chosen queries $\mathbf{q} = (q_1, q_2, q_3, \dots, q_o)$ and receives from \mathcal{C} a set of tokens $\mathbf{t} = (\text{tk}_1, \text{tk}_2, \text{tk}_3, \dots, \text{tk}_o)$.
- Finally \mathcal{A} outputs $v = (\text{SES}, \mathbf{t})$.

$\mathbf{Ideal}_{\mathcal{A}, \mathcal{S}(\lambda)}^{S^3E}$ game:

- \mathcal{A} outputs a string stream S .
- The simulator \mathcal{S} through the leakage \mathcal{L} generates SES and forwards it to \mathcal{A} .
- \mathcal{A} issues a polynomial number of queries $\mathbf{q} = (q_1, q_2, q_3, \dots, q_o)$. \mathcal{S} replies to each of the queries through the leakage function \mathcal{L} with $\mathbf{t} = (\text{tk}_1, \text{tk}_2, \text{tk}_3, \dots, \text{tk}_o)$.
- Finally \mathcal{A} outputs $v = (\text{SES}, \mathbf{t})$.

Definition 3. A S^3E scheme is adaptively \mathcal{L} -semantically secure against a probabilistic polynomial time adversary \mathcal{A} if there is exists a polynomial Simulator \mathcal{S} such that for all polynomial time distinguishers \mathcal{D} :

$$|\Pr[\mathcal{D}(v) = 1 : v \leftarrow \mathbf{Real}_{\mathcal{A}(\lambda)}^{S^3E}] - \Pr[\mathcal{D}(v) = 1 : v \leftarrow \mathbf{Ideal}_{\mathcal{A}, \mathcal{S}(\lambda)}^{S^3E}]| \leq \text{neg}(\lambda)$$

3 Related work

The ORAM paradigm [21] enables a user to remotely search for encrypted data, without leaking the search or the access pattern. The trade off comes with a bandwidth and communication burden. In [21] the bandwidth overhead is polylogarithmic, which has been reduced down to logarithmic in subsequent work [14, 34, 37, 39, 40]. However, in order to provide a real practical real world remote search protocol on encrypted data some leakages are allowed: the *search* and *access* pattern. The formalization of these patterns and efficient schemes have been presented in the literature under the **Symmetric Searchable Encryption** (SSE) framework [6–8, 12, 36].

With SSE a user encrypts data and index separately. It uploads both to an untrusted cloud and later on can search efficiently file identifiers with specific single keywords or an expressive boolean function over keywords, without the cloud learning anything about the files or the keywords. This comes at a security cost of leaking the *search* and *access pattern*. Following the approach of SSE, we can design substring searchable symmetric schemes as follows. The user builds an index which maps substrings to positions, encrypts the index and uploads it to the cloud. Later on, the user computes a token for the specific substring and the cloud tries to find a match in the index. If a match occurs the cloud returns the encrypted positions for this token, which correspond to a substring. However, this approach has increased bandwidth cost $O(n^2)$, since the cloud has to keep track of all the possible substrings.

Tailored substring searchable encryptions schemes have been proposed in the literature [9], [16], [15]. Here we present a detailed analysis of the state of the art in substring searchable encryption protocols. Chase *et al.* [9] leverage the auxiliary data structure of the suffix tree. A suffix tree is a compressed suffix trie, can be computed in time $O(n)$ and allows for substring search in $O(m)$ time on a substring of size m . Its amortized storage cost is $O(n)$ which hinders a big constant factor, which can go up to 20 [1, 5, 22, 30]. In [15] the authors extended the efficient SSE scheme for boolean queries from [6] in order to support substring matching. The idea is to build an index of overlapping k -grams, to prepend its relevant position and encrypt it. When a user needs to perform a substring query, the cloud performs a conjunctive keyword search for all the k -grams of the substring and returns the position. The disadvantage of the scheme comes at the need of storing all the overlapping k -grams at the cloud, which will represent substrings.

In [16] the authors follow a different approach. Instead of taking the *index-then-encrypt* approach with fast symmetric cryptographic primitives, they modify the subset sum problem, which is used to build public key encryption schemes, in a means such that the cloud can solve it. This contradicts the security definition of subset sum problem, which asks for impossibility of an adversary to find a solution to a specific instance. More specifically the user uploads a special instance \mathbf{T} of a subset sum problem such that given a trapdoor R_i associated with a substring, the solution can be solved in time $O(m)$ by the cloud; with the special property that the integers which sum to \mathbf{T} , parametrized by R_i are the positions of the substring in the string. This technique hides also the search pattern but comes at the cost of fixed size substrings, that must be defined in the beginning of the protocol. Moreover the substring should be substantially small with respect to the big stream. Our solution in contrast allows variable size of substring of any size.

Recently, Blass and Moataz [33] strengthen the security requirements by hiding the search and access patterns, following the ORAM approach. By leveraging the Path ORAM technique and the suffix array construction for substring queries, the authors manage to reduce the bandwidth, with a binary recursive tree above the position map. Each node in the tree represents a Path ORAM of the binary search tree for the suffix array. However, in order the cloud to be able to perform an oblivious binary search has to keep track of all the suffixes, which blows up the storage cost for the server. Furthermore the need for storing the suffixes cancels out the suffix array storage advantage over suffix tree. Finally, due to the Path ORAM technique the user has to store a state logarithmic on the length of the string—for the position map. The extra security guarantees of the tailored ORAM scheme do not allow for efficient storage cost both at the client and the cloud size, which is the goal for our work.

Papadopoulos *et al.* [35] addressed the problem of authenticating substring queries without privacy and various work for pattern matching adopts the two party computation model [13, 18, 23, 24] in which one party holds the data stream and a client the pattern. The model differs from the substring searchable symmetric encryption, since in the latter one client holds both the pattern and the stream and uploads an index of the stream to an untrusted party.

4 Idea

In order to reduce the storage cost for our Substring Searchable Symmetric encryption scheme (S^3E) we first substitute the storage expensive suffix tree of the state of art work in [9] with a suffix array SA . A suffix array for a string S of size n constitutes of an integer array of size n , which has at each position a pointer to the start of the matched suffix $T[1..m]$ in the string S . SA is lexicographically sorted with respect to all the possible suffixes and can be computed in linear time on the size of the string S . In order to look for the position of a substring, a binary search in SA is performed, which is used as an index to the original string. Thus, the running time for a substring search is $O(m + \log n)$. Let us now consider a concrete example to uncover its details. Suppose $S=lalakis$. The algorithm for the suffix array proceeds as follows:

1. Compute all the suffixes starting from the right-most position: s, is, kis, akis, lakis, alakis, lalakis.
2. Lexicographically sort the suffixes: akis, alakis, is, kis, lakis, lalakis, s.
3. Find the position in S of each suffix from step 2 and store them in an array $SA = [4, 2, 6, 5, 3, 1]$
4. Output SA .

However plugging the SA for a substring searchable symmetric encryption scheme raises some difficulties. We assume that the suffix array is encrypted under a secret key of the user. In order the cloud to retrieve the right encrypted index position from SA ought to run a binary search obviously without learning the underlying string S , query substring T , or any of the suffixes. A solution to the problem is to use the technique presented by Gentry *et al.* [19], which allows for a single ORAM query in order to perform a binary search over encrypted data. However, in order to adapt this approach it is required from the server apart from the encrypted suffix array, to store the tree of the encrypted data, which would be an extra burden for its storage complexity, plus there should be one extra round of communication due to the ORAM protocol in order to rebuild the specific path.

In this paper we are taking a different approach, which achieves storage efficiency. The problem arises from the fact that the server does not hold the original string in order to perform a binary search, which is indexed by the suffix array SA . We take advantage of a *self-indexed* data structure called hereafter FM index [17]. Namely, from FM index the untrusted cloud can answer substring queries by leveraging the suffix array SA , without the need for an ORAM query. The neat property of the FM index is that it can reconstruct the original string S with some extra auxiliary data structures, thanks to its instantiation from the BWT compression algorithm [5]. For the reconstruction it employees the LF mapping technique, thus there is no need to store the encrypted stream S . The FM index can be derived from SA , as such its computational overhead is almost for free, after the computation of the suffix array. We describe the core building blocks of the FM index in what it follows.

4.1 Pattern matching

In this section we describe the compressed index FM, that will be used for the construction of our Secure Pattern Matching (S^3E) protocol. The design lies heavily on the BWT transformation for compression of bit-strings and on a special LF mapping for the reconstruction of the original string from BWT. The BWT, along with the LF mapping technique and some auxiliary information are the basic blocks of the compressed index for substring queries.

BWT Transformation The Burrows-Wheeler Transformation (BWT) transforms a stream of data by leveraging the entropy of each character. In a nutshell the data stream S is transformed to an encoding W such that compression algorithms provide high rate of compression. However for the construction of S^3E we will only need a compressed version of the intermediate steps and not the final string W . For ease of completeness we show the steps to transform an original stream S to W with BWT in algorithm 1. First the algorithm appends the terminating symbol $\$$ to the input string S . Then it builds the matrix W by permuting the symbol $\$$. At each iteration the permutation is appended as a new row to the matrix W . Finally the rows of W are sorted lexicographically in an ascending way. A real world example is shown in figure 2 for string $S = lalakis$. The upper table of the figure shows the permutations and the final sorted matrix W is shown at the bottom matrix. The transformation is the first step towards compression with the LF mapping that is shown below.

Algorithm 1: BWT transformation

Input: String S
Output: $BWT(S) = W$
 $l = \text{length}(S) + 1$;
 $S.append(\$)$;
 $i = 0$;
while $i < l$ **do**
 $r_i = \text{rotate}(s, \$)$ // The rotate algorithm permutes the characters of the original string and returns the permuted string;
 $W.addrow(r_i)$ // It adds the permuted row from the previous step to the matrix W ;
 $i++$;
end
return Sorted. W ;

l	a	l	a	k	i	s	\$	\$	l	a	l	a	k	i	s
a	l	a	k	i	s	\$	l	a	k	i	s	\$	l	a	l
l	a	k	i	s	\$	l	a	a	l	a	k	i	s	\$	l
a	k	i	s	\$	l	a	l	i	s	\$	l	a	l	a	k
k	i	s	\$	l	a	l	a	k	i	s	\$	l	a	l	a
i	s	\$	l	a	l	a	k	l	a	k	i	s	\$	l	a
s	\$	l	a	l	a	k	i	l	a	l	a	k	i	s	\$
\$	l	a	l	a	k	i	s	\$	l	a	l	a	k	i	s

Fig. 2: BWT Transformation. The upper table shows the cyclic permutations for the string $S=lalakis$. As a first step for the construction of the BWT matrix the end symbol $\$$ is appended at the end of S , which lexicographically precedes all alphabetical symbols. Then at each row a permutation of $\$$ around the string is shown. The result is matrix $BWM[\text{length}(S)+1][\text{length}(S)+1]$. As a second step the rows are sorted lexicographically at the second bottom table of the figure starting from the first character of each string. The light gray first column shows the order of the characters after sorting and the last column is the result of the transformation $BWT(S)=sllkaa\$s$.

LF Mapping The LF Mapping technique takes the first F and last L columns from the BWT transformation and through an iterative process as described in algorithm 2 reconstructs the original string S . Starting from the first elements of each column from F and L , the algorithm employees L as an index to the F column. Each time the element of the L column is appended to a LIFO stack. The value at the current position will be used as an index for the F column for the next loop. An example is presented in figure 3. At the first iteration the pointer indicates the first position in both columns F, L . For the next iteration the L character 's' indicates the index for the first column F , which can be found at its last position with $F[7]=s$. The current character at the L column is appended to a stack D . For a next iteration the current character at the L column indicates the next index for the F column. The character i is pushed to the stack D . The procedure halts when the position at L is $\$$. Then the algorithm pops all elements from D and the initial string S is fetched.

	F	L	F	L	F	L	F	L	F	L	F	L	F	L	F	L
	\$	s	\$	s	\$	s	\$	s	\$	s	\$	s	\$	s	\$	s
a	l	a	l	a	l	a	l	a	l	a	l	a	l	a	l	a
a	l	a	l	a	l	a	l	a	l	a	l	a	l	a	l	a
i	k	i	k	i	k	i	k	i	k	i	k	i	k	i	k	i
k	a	k	a	k	a	k	a	k	a	k	a	k	a	k	a	k
l	a	l	a	l	a	l	a	l	a	l	a	l	a	l	a	l
l	\$	l	\$	l	\$	l	\$	l	\$	l	\$	l	\$	l	\$	l
s	i	s	i	s	i	s	i	s	i	s	i	s	i	s	i	s

Fig. 3: The LF mapping process is used to reconstruct the original string S from the transformed one after applying the BWT operation. Starting with the $\$$ sign from the F column, the mapping progressively reconstructs the entire string S . The last column L is used as a "ladder step" to find the next i th index in the F column, which in turn maps to the i th entry in the L column. The entire procedure halts when $L[i]=\$$

Algorithm 2: LF Mapping

Input: First (F), Last column (L) from BWT

Output: S

$D=0$ // Initialize the stack D;

$l=length.(F)$ // the length of F equals the length of L;

$i=0$;

while $L[i] \neq \$$ **do**

$D.push(L[i]);$

$i=find.F[L[i]]$ // $find.[\]$ denotes the index number in array $[\]$ that the element is. For instance

$find.F['s']=7$;

while $D \neq \backslash '0'$ **do**

$S=S+D.pop$;

return S ;

FM Index Suffix array vanilla construction has $O(n^2 \log n)$ asymptotic computational cost. This stems from the fact that we need to first sort the n suffixes by performing $O(n \log n)$ comparisons and each comparison has cost n . Linear time algorithms have been achieved by first constructing a suffix tree and then traversing with a depth first edges with lexicographical order. However our goal is to be storage efficient, meaning we want to eliminate the storage cost of a suffix tree which practically approximates a constant factor of $20n$ [1, 5, 22, 30]. We pick up the *skew* algorithm [26] which is a *divide and conquer* based algorithm and achieves linear time construction. The approach of the *skew* algorithm is to recursively divide the suffixes in three groups depending on the position pos of all suffixes: $pos \bmod j, j \in \{1, 2, 3\}$ and then merge the result.

The FM consists of three column arrays. The first one is the F column from the LF mapping, the second one is the L column which corresponds to the BWT(S) and the last one is the position of each substring to S , which is obtained from the suffix array SA. Notice that the F column is equivalent with the same range first letter suffixes of the corresponding SA. $L = BWT(S)$ can be computed with the formula $BWT(S)[i] = S[SA[i] - 1]$ from the suffix array. Furthermore for the traversal of the LF mapping the unique ranking of each character in each F, L needs to be stored in r_F, r_L accordingly. Finally $FM = \{F[i], L[i], r_F[i], r_L[i], SA[i]\}_{i=1}^n$

5 Protocol

In this section we progressively show the design of the Substring Searchable Symmetric encryption scheme S^3E in order to facilitate its design principles. We give a description of our first approach. We identify its subtleties and we explain why the scheme leaks more than the allowed leakage \mathcal{L} as defined in Section 2.

First we give a naive construction which is vulnerable to **offline frequency attacks**; meaning the adversary is able to perform a frequency attack based on the encrypted index without obtaining any search token by the client. We then, mitigate this attack and we allow only **online frequency attacks**, in which the adversary is able to learn the frequency of the characters based on the search query. However, after many queries the honest-but-curious cloud can eventually decrypt the index which contains the original data stream. We furthermore change the scheme in order to be resistant to online frequency attacks.

Despite the lack of frequency attacks, we usher a novel attack mechanism owing to the inherent construction of the FM index. First we identify **offline-difference-attack**, which permits the untrusted cloud to eventually recover all encrypted positions of the suffix array SA. To partially circumvent this attack we alter our first approach in order to transform **offline-difference-attack** to **online-difference-attack**, in which the cloud needs n queries to be issued by the client in order to recover the encrypted SA. Finally we present the full details of the protocol which curbs any difference attack.

5.1 First Approach

For our S^3E protocol we are using a hash map $LLSet$, which is a set of tuples $\langle k, v \rangle$, with keys k to access values v . $LLSet$ values entail the addresses for a set of linked lists LL_k and an array FM. For the security

of the scheme the user employs lightweight cryptographic primitives: a pseudorandom function $F(\cdot)$, a pseudorandom permutation $\Pi(\cdot)$ respectively and a symmetric encryption scheme $\text{SKE} = \{\text{Gen}, \text{Enc}, \text{Dec}\}$.

During the **KeyGen** phase of the protocol user chooses the secret key $\mathbf{k} = (k_f, r, k_\pi, k_e)$. Keys k_f, r is for the pseudorandom function $F(\cdot)$, k_π for the pseudorandom permutation $\Pi(\cdot)$ and k_e for the symmetric encryption scheme $\text{SKE} = \{\text{Gen}, \text{Enc}, \text{Dec}\}$. During the **PreProcess** step the user computes a hash table of linked lists LLSet , where each position $\text{LLSet}[F_{k_f}(c)]$ maps to the linked list LL_c . The number of linked lists equals the number of distinct elements c , denoted as w in the data stream S , where symbol c comes from an alphabet Σ . The hash table is used to fetch all the positions of a character in the stream S from the linked lists LL_n . Each linked list LL_c stores information concerning the retrieval of the position of c from S . More specifically each node in the list stores the following tuple: $\langle \text{nptr}, \text{addr} \rangle$, nptr is a pointer to the next node of the current list and addr is the address of the element c in the **FM** index. The **FM** index is a three dimensional array which keeps track of the **F, L** columns and the encrypted **SA** suffix array, which contains only encrypted positions of substrings and not all the suffixes as in [33]. The untrusted cloud, thanks to the **LF** mapping and the **FM** index computation does not need to store all the suffixes of a stream S (cf. figure 4). To recap, the user computes the suffix-array **SA** and the **F, L** columns through the **BWT** transformation. The set of tuples $\langle F_j, L_j, \text{SA}[j] \rangle_{j=1}^w$ are encrypted using a secure PRF $F_{k_f}(\cdot)$ and a symmetric encryption algorithm $\text{SKE} = \{\text{Gen}, \text{Enc}, \text{Dec}\}$ as follows: Using the PRF $F_{k_f}(\cdot)$ user computes $t_j = \langle F_{k_f}(c_{F_j}) \oplus F_{k_f}(r_{F_j}), F_{k_f}(c_{L_{c_j}}) \oplus F_{k_f}(r_{F_{c_j}}) || F_{k_f}(r_{L_{c_j}}), \text{Enc}(\text{pos}_j) \rangle_{j=1}^k$, where c_{F_j} denotes the character at position number j of the **F** column from the **BWT** transformation, c_{L_j} the j th character at the **L** column and pos_j is the index number in the suffix array which dictates the position in stream S which is indexed in position j in the suffix array **SA**.

The crux of the design is on how to allow fast indexing through a hash table, which means that there should be unique keys derived from the string with repetitive characters. We employ the ranking information r_c of each character along with the character itself. This coupling makes a unique bit-stream which can be given as input to the PRF $F_{k_f}(\cdot)$ and serves as an index to the **FM** design. However, when a user is looking for a substring, she does not know the ranking of each character in the substring $T[1..m]$. We mitigate this deficiency by building an linked list LL_c for each character. Each node in the linked list maps to the address of the character c in the **FM** index, and has a pointer to the next same character node. The first node of each linked list is stored in the LLSet hash table. In order to prevent frequency attacks we encrypt each key $F_{k_f}(c)$ in the LLSet hash map with another key r as follows: $F_{k_f}(c) \oplus F_r(c)$. Thus the cloud cannot perform a frequency attack offline without observing any token. The key of the hash map LLSet at $F_{k_f}(c) \oplus F_r(c)$ maps to the first element of the linked list LL_c , which is encrypted as $\langle \text{nptr}, \text{addr} \rangle \oplus F_r(c)$. As such, the frequency of each character before a search query is hidden. The second difficulty comes when the cloud tries to traverse the **FM** index through the **LF** mapping technique. The encrypted **FM** index contains unique digests of characters, while the cloud should identify matches from the token $\text{tk}_{T,S}$, that encodes repetitive characters deterministically. In order to allow the cloud traverse the encrypted **FM** index, we encrypt the **FM** as a key value hash table where the key consists of $F_{k_f}(c_{F_j}) \oplus F_{k_f}(r_{F_j})$ and the value is $F_{k_f}(c_{L_{c_j}}) \oplus F_{k_f}(r_{F_{c_j}}) || F_{k_f}(r_{L_{c_j}}), \text{Enc}(\text{pos}_j)$. Conceptually the key $F_{k_f}(c_{F_j}) \oplus F_{k_f}(r_{F_j})$ encodes the **F** column, $F_{k_f}(c_{L_{c_j}}) \oplus F_{k_f}(r_{F_{c_j}}) || F_{k_f}(r_{L_{c_j}})$ encrypts the **L** column and $\text{Enc}(\text{pos}_j)$ the suffix array. Intuitively $F_{k_f}(r_{F_{c_j}})$ at the **L** column acts as a key for the $F_{k_f}(c_{L_{c_j}})$ value. The cloud decrypts it only upon receipt of a substring token. Thus we prevent offline attacks in which the cloud can learn the frequency of the characters without any token from the user. Notice that the cloud after observing a token, can learn the encryption of the rankings which is unique for all characters $F_{k_f}(r_{F_j})$ and will try to xor it with all elements in the **F** column in order to perform a frequency analysis. We can mitigate this attack by encrypting unique rankings per character as follows: $F_{k_f}(r_{F_j} || c_j)$. Thus the user encrypts the **FM** index as: $t_j = \langle F_{k_f}(c_{F_j}) \oplus F_{k_f}(r_{F_j} || c_{F_j}), F_{k_f}(c_{L_{c_j}}) \oplus F_{k_f}(r_{F_{c_j}} || c_{F_j}) || F_{k_f}(r_{L_{c_j}} || c_{L_j}), \text{Enc}(\text{pos}_j) \rangle_{j=1}^k$. Finally the user permutes all the tuples with a secure permutation: $\Pi_{k_\pi}(t_j)$.

Even though we circumvent **offline frequency attacks**, by encrypting appropriately the set of linked lists, **online frequency attacks** cannot be accepted as a security leakage in a real world scenario, since the cloud can decrypt queries and the encrypted index through frequency analysis, by observing search queries. In order to blur frequency analysis on the encrypted index from substring search queries we pad each list LL_c with elements up to the maximum size list LL_c with pointers to the **FM** index at consecutive positions starting from the end. This will negatively affect the correctness of the protocol since the cloud now will erroneously reply with matching queries for non existent character streams. We solve this problem by extending the **FM** index with encrypted tuples of same alphabet characters, which

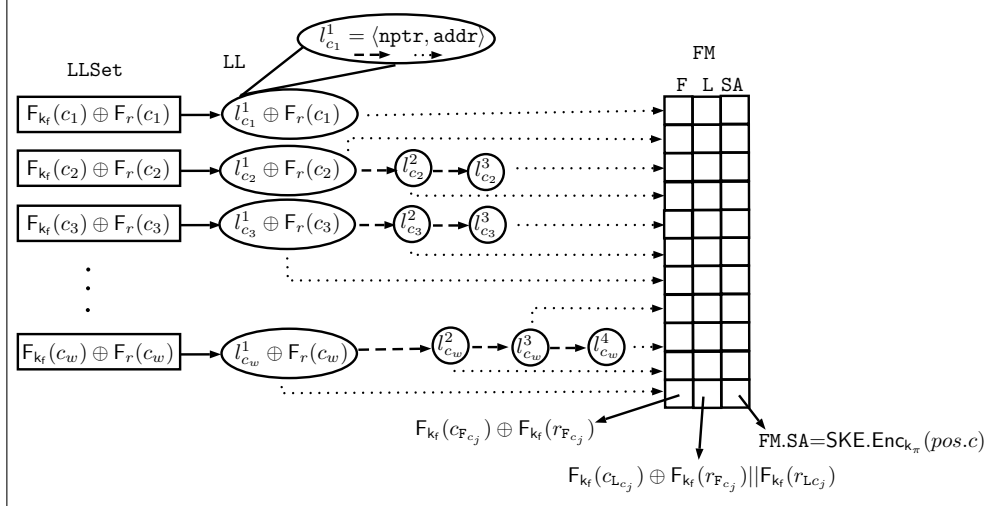


Fig. 4: The encrypted FM index construction.

are indistinguishable from the valid data stream. However, this will raise false positive responses. The client can verify the correctness and discard false positive responses by checking the returned suffix array position pos : if $\text{pos} > n$ then the substring occurs from lazy padding, otherwise it is a valid position.

Notice that flushing all the elements of all lists in a big array at (pseudo) random positions as in [12] will not mitigate the problem. After observing a substring token the cloud learns its frequency and can perform a frequency attack. This leakage is not devastating in [12] as the revealed information is the cardinality of file identifiers for an encrypted document and not the frequency of a character as it is in a substring searchable symmetric encryption scheme. Thus, padding the lists is imperative in order to tackle **online frequency attacks**.

During the Search phase on a substring query $\text{tk}_{T,S} = F_{k_f}(T[1..m]) = F_{k_f}(T[1]), F_{k_f}(T[2]), \dots, F_{k_f}(T[m]), F_r(T[m])$ the cloud proceeds as follows: From the LLSet hash table it looks for the value with key $F_{k_f}(T[m]) \oplus F_r(T[m])$. This value maps to a linked list LL_c , in which each node maps to the encrypted FM index $t_j = \underbrace{\langle F_{k_f}(c_{F_j}) \oplus F_{k_f}(r_{F_j} || c_{F_j}) \rangle}_F, \underbrace{\langle F_{k_f}(c_{L_{c_j}}) \oplus F_{k_f}(r_{F_{c_j}} || c_{F_j}) || F_{k_f}(r_{L_{c_j}} || c_{L_j}) \rangle}_L$,

$\underbrace{\text{Enc}(\text{pos}_j)}_{\text{SA}} \Big|_{j=1}^k$. In order to decrypt the first element of the linked list the cloud uses $F_r(T[m])$ as a key

to decrypt $\langle \text{nptr}, \text{addr} \rangle \oplus F_r(c)$, in order to learn $\langle \text{nptr}, \text{addr} \rangle$. The cloud uses $F_{k_f}(T[m])$ and applies a xor operation on the F column at the ranges that it retrieved from the linked list of the c_m character LL_c and learns $F_{k_f}(r_{F_j} || c_{F_j})$. Afterwards it uses $F_{k_f}(r_{F_j} || c_{F_j})$ as a key to decrypt the first part of the L column element $F_{k_f}(c_{L_{c_j}}) \oplus F_{k_f}(r_{F_{c_j}} || c_{F_j})$ and reveals $F_{k_f}(c_{L_{c_j}})$. It then fetches the encrypted L column as $k = F_{k_f}(c_{L_j}), r = F_{k_f}(r_{L_j} || c_{L_j})$ in which $F_{k_f}(c_{L_j}) = F_{k_f}(T[m-1])$ and for all nodes from the linked list computes $k \oplus r$, which is used as a key for the F column. The procedure terminates when the processed substring character is the first one $F_{k_f}(T[1])$. At this point the cloud returns to the user all the encrypted $\text{Enc}(\text{pos}_j)$ for the substrings. The user decrypts accepts the result as long as the decrypted position is in the range of the size of original stream without padding.

Offline-difference-attack(ofda). The above design is vulnerable to **ofda**. Namely, from the per-character one way function F_{k_f} evaluation of the substring query: $\text{tk}_{T,S} = C[1], C[2], \dots, C[m] \leftarrow F_{k_f}(T[1..m])$ and the LF mapping the protocol leaks to the cloud in cleartext the exact differences of the positions of two encrypted substrings in the stream S . More specifically, the number of iterations in the LF mapping traversal, reveals how many positions two substrings they differ, as long as there is match in S . Eventually, an untrusted cloud can decrypt the entire encrypted SA array, which contains encrypted positions of all substrings in S , since it knows its addresses. Notice, that the attack is characterized as offline, in the sense that the cloud does not interact with the user and can launch the attack prior to the issue of user's substring queries.

Online-difference-attack(onda). To mitigate the aforementioned leakage, we alter the protocol and the construction of the index. The user instead of coupling its encrypted cell of the SA array with the same range cells of the encrypted F and L columns, it permutes separately the SA array and stores only the addresses of the cells in the tuple $(F_{k_f}(c_{F_j}) \oplus F_{k_f}(r_{F_j} || c_{F_j}), F_{k_f}(c_{L_{c_j}}) \oplus F_{k_f}(r_{F_{c_j}} || c_{F_j}) || F_{k_f}(r_{L_{c_j}} || c_{L_j}), \text{Enc}(\text{addr}(\text{pos}_j)))_{j=1}^k$. During the Search phase the cloud returns the encrypted address of the possible match(es). Finally the user decrypts, fetches the encrypted position(s) with the aid of the cloud, decrypts and learns the position(s) in S . It verifies the validity of the result depending on the range of the position with respect to the size n of the stream: $\text{pos} > n$ denotes a false positive and non-valid match. Thanks to the encrypted addresses of possible matched position(s) in S and the extra round of communication between the cloud and the user, the new protocol effaces **ofda**.

However, the cloud can still identify the differences in positions in S of two substring queries, as long as the user issues the queries. Eventually after many queries the cloud can fully decrypt the encrypted SA, since it contains permuted unique encryptions of $1..n$. In order to expunge the online-difference-attack we add a private information retrieval protocol at the last round, in which the user retrieves the correct encrypted position, without the cloud learning its address, meaning both **ofda** and **onda** cannot threaten the scheme.

Attacks	Mitigation
Offline Frequency	Extra character encryption with key r
Online Frequency	List Padding
Offline Difference	Permute and Encrypt Suffix array address
Online Difference	PIR

Table 1: Attacks and mitigation methods.

5.2 S³E Description

We are now ready to give the full details of our substring searchable symmetric encryption protocol S³E in figures 5 and 6.

5.3 Security Analysis

We illustrate the security of the scheme pertaining to definition 3. More specifically we show the existence of a simulator \mathcal{S} who has access to the leakage function $\mathcal{L} = (\mathcal{L}_1, \mathcal{L}_2, \mathcal{L}_3)$, and produces indistinguishable views to an adversary \mathcal{A} . Conceptually the proof demonstrates that an adversary \mathcal{A} , who can be a dishonest cloud cannot learn partially more information from what it can be leaked in an ideal work without malicious behaviors. Albeit the acceptable introduction of a leakage which is permitted in the proof, the technique has been broadly adopted in the queryable encryption schemes. Following the same trajectory, we are not questioning the implications of the leakage in the security of the scheme as long as they have been rigorously defined. Further analysis is needed however in order to illustrate in clearance the leakages in substring symmetric encryption schemes owing to the assembly of intricate abstract data types.

Theorem 1. *Let $F_{k_f}, \Pi_{k_\pi}, \text{SKE} = \{\text{Gen}, \text{Enc}, \text{Dec}\}, \text{PIR}$ be a pseudorandom function, a pseudorandom permutation, a semantically secure symmetric encryption scheme and a secure PIR scheme respectively, then our substring searchable symmetric encryption scheme S³E is adaptively \mathcal{L} -semantically secure.*

The proof is omitted to Appendix B section.

5.4 Comparison

We perform a comparison of our S³E with existing solutions (cf. table 2). We analyzed the search running time in asymptotic complexity, index space requirements both in the plaintext and in the ciphertext space, query size, variable length capability, rounds of communication and search leakage. Since our scheme competes mostly with [9] here we further elaborate its cost analysis from table 2.

- $k \leftarrow \text{KeyGen}(1^\lambda)$: This algorithm run by the user takes as input the security parameter 1^λ and generates random keys $k = (k_f, r, k_\pi, k_e)$ for a PRF $F_{k_f, r} : \{0, 1\}^\lambda \times \{0, 1\}^\nu \rightarrow \{0, 1\}^\mu$, a PRP $\Pi_{k_\pi} : \{0, 1\}^\lambda \times \{0, 1\}^\nu \rightarrow \{0, 1\}^\nu$ and a symmetric encryption algorithm $\text{SKE} = \{\text{Gen}, \text{Enc}, \text{Dec}\}$. Finally it outputs k to the user. For the generation of the keys we assume a source of randomness \mathcal{R} and a pseudorandom generator G seeded with $s_f \xleftarrow{\mathcal{S}} \mathcal{R}, s_r \xleftarrow{\mathcal{S}} \mathcal{R}, s_\pi \xleftarrow{\mathcal{S}} \mathcal{R}, s_e \xleftarrow{\mathcal{S}} \mathcal{R} : (k_f, r, k_\pi, k_e) \leftarrow G(s_f), G(s_r)G(s_\pi), G(s_e)$.
-
- $\text{SES} \leftarrow \text{PreProcess}(k, S)$: User owns a stream S , which contains characters $c \in \Sigma$. S has n characters in total and k distinct elements. User:
 1. Computes the suffix array SA and the F, L columns and stores it as the FM index: $\text{FM} = \text{F} \parallel \text{L}, \text{SA}$.
 2. Encrypts its element of SA array with $\text{SKE.Enc}(k_e, \text{SA}[i]), 1 \leq i \leq n$ and encrypts the addresses of each element $\text{SA}[i]$ at a new array, ASA using $\text{Enc}(\text{ASA}[i]), 1 \leq i \leq n$.
 3. Applies the PRF to each element of F as follows: $\text{F}[i] = F_{k_f}(c_{\text{F}_i}) \oplus F_{k_f}(r_{\text{F}_i} \parallel c_{\text{F}_i})$.
 4. Computes $\text{L}[i] = F_{k_f}(c_{\text{L}_{c_i}}) \oplus F_{k_f}(r_{\text{F}_{c_i}} \parallel c_{\text{F}_i}) \parallel F_{k_f}(r_{\text{L}_{c_i}} \parallel c_{\text{L}_i})$.
 5. Applies a pseudorandom permutation Π_{k_π} to the tuples $t_i = \langle F_{k_f}(c_{\text{F}_i}) \oplus F_{k_f}(r_{\text{F}_i} \parallel c_{\text{F}_i}), F_{k_f}(c_{\text{L}_{c_i}}) \oplus F_{k_f}(r_{\text{F}_{c_i}} \parallel c_{\text{F}_i}) \parallel F_{k_f}(r_{\text{L}_{c_i}} \parallel c_{\text{L}_i}), \text{Enc}(\text{ASA}[j]) \rangle_{j=1}^n : \Pi_{k_\pi}(t_i) = \{t'_i\}_{j=1}^n = \text{FM}'$.
 6. For every distinct character in $\text{F}[i] = F_{k_f}(c_{\text{F}_i}) \oplus F_{k_f}(r_{\text{F}_i} \parallel c_{\text{F}_i})$ the user builds a linked list LL_c and each node stores $\text{LL}_c.\text{nptr}$ for the next node of the list and $\text{LL}_c.\text{addr}$ which points to the tuple t_i with a matching $F_{k_f}(c_{\text{F}_i})$. Finally it encrypts the first element of each linked list LL_c with $F_r(c_i) : \langle \text{nptr}, \text{addr} \rangle \oplus F_r(c_i)$.
 7. Let max_c be the cardinality of most frequent character and f_{c_i} the frequency of character c_i . User chooses dummy characters $\{dc\}_{i=1}^{\sum_{j=1}^k \text{max}_c - f_{c_j}}$ that constitute a dummy stream. Computes the FM index of it, encrypts it as in steps 2,3 and 4. It appends the new dummy index DFM to the existing one $\text{FM}' = \text{FM}' + \text{DFM}$.
 8. The head pointers of the collections of all linked lists are stored in a hash table LLSet with key $k = F_{k_f}(c_i) \oplus F_r(c_i)$ and value v a pointer to the head of the list LL_c , which stores information about the $F_{k_f}(c_{c_i})$ character, meaning all its positions to the encrypted FM index.
 9. Finally outputs $\text{SES} = (\text{LLSet}, \text{LL}_c, \text{FM}')$ and keeps only the keys $k = (k_f, r, k_\pi, k_e)$.

Fig. 5: The KeyGen and PreProcess algorithms of the S^3E protocol.

Thanks to the usage of encrypted dictionary the cost of searching a m length string is $O(m + k)$, where k denotes the number of occurrences. However due to the padding technique, which mitigates *online frequency attacks* and the PIR protocol during the last round the search cost is increased to $O(m + k + \sum_{j=1}^k \text{max}_c - f_{c_j} + p(n))$, where max_c is the most frequent character and f_{c_j} the frequency of character c_j and $p(n)$ stands for the PIR search cost. Assuming a parallel implementation of PIR on GPUs the cost can be parallelized with multiple running threads thus reducing a lot the constant factor from $O(n)$ to $p(n)$ with $p(n) \ll O(n)$ initial cost of PIR.

For the index space complexity we analyzed the space requirement in the plaintext space and in the ciphertext space. For the plaintext space analysis we assume that a pointer or integer requires 4 bytes. Recall that a suffix tree has n leaves, at most $n - 1$ internal nodes and at most $2n - 2$ edges. Thus, for a naive suffix tree implementation we need 2 pointers for each leaf: one for the parent node and one for its position to the original stream, resulting in $8n$ bytes. Four pointers for each internal node: one for the parent node, one for each leftmost child, one for the right sibling and one pointer for the suffix link, which reduces the search time during a substring query. The total storage cost for the internal nodes is $4 * 4n = 16n$. For each edge, suffix trees allocate one pointer for the beginning position of the substring in the stream and one for the end position of the substring in the stream increasing the space cost to

- $\text{tk}_{T,S} \leftarrow \text{SrchToken}(k, T[1\dots m])$: This algorithm takes as input the secret substring search key k , a string $T[1\dots m]$ and outputs a trapdoor to search for the string T on data stream S , through SES . User with his secret PRF key k_f computes $\text{tk}_{T,S} = C[1], C[2], \dots, C[m] \leftarrow F_{k_f}(T[1\dots m]), F_r(T[m])$ and forwards $\text{tk}_{T,S}$ to the cloud.
-
- $(s, \perp) \leftarrow \text{Search}(\text{tk}_{T,S}, \text{SES})$: The cloud parses the token query $\text{tk}_{T,S} = C[1], C[2], \dots, C[m]$ and searches the position in S from the encrypted index SES as follows:
 1. $u = \text{find}(\text{LLSet}, C[m])$ // find in dictionary LLSet the value u with key $C[m]$. u is a pointer to the head of a list LL , which stores pointers to all characters $T[m]$ in S
 2. **if** $u == \perp$ **return** \perp
 3. **while** $u \neq \perp$ **do**
 - $K = K \cup u.\text{addr}$ // traverse the list and store in the set K the addresses of the characters.
 - $u = u.\text{nptr}$
 4. **for** $p = m - 1; p > 1; p = p - 2$
 - for** $i = 1; i < \text{size}(K); i ++$
 - if** $\text{SES.L}[K[i]]^{(1)} == C[p - 1]$ // Store in the set KEYS only the elements from the F column, whose associated L element equals the next character from C in a backward order. $\text{SES.L}[K[i]]^{(1)}$ maps to $F_{k_f}(c_{L_{c_i}}) \oplus F_{k_f}(r_{F_{c_i}} || c_{F_j})$ and $\text{SES.L}[K[i]]^{(2)}$ to $F_{k_f}(r_{L_i})$.
 - $\text{KEYS} = \text{KEYS} \cup \text{SES.L}[K[i]]$
 - else** $K = K - K[i]$ // Remove all the non matched elements from the key set K .
 - if** $K == \perp$ **return** \perp
 - for** $i = 1; i < \text{size}(\text{KEYS}); i ++$
 - $r_{F_{c_i}} = C[p] \oplus \text{KEYS}^{(1)}[i]$
 - $z = r_{F_{c_i}} \oplus \text{KEYS}^{(2)}[i]$ // Compute the key from the L column as $\text{KEYS}^{(1)}[i] \oplus \text{KEYS}^{(2)}[i]$, which corresponds to $F_{k_f}(c_{F_i}) \oplus F_{k_f}(r_{F_i} || c_{F_j})$ in the F column of the SES object.
 - if** $\text{SES.F}[z] \neq \perp$
 - continue**
 - else** $\text{KEYS} = \text{KEYS} - \text{KEYS}[i]$ // Remove all the non matched elements from the key set KEYS .
 - $K = \text{KEYS}$
 - 5. **if** $K == \perp$ **return** \perp
 - 6. Cloud sends to the user $\text{SES.ASA}[K]$
 - 7. Client decrypts with k_e all the elements of $\text{SES.ASA}[K]$, that correspond to addresses that contain encrypted positions of the requested substring in S : $\mathbf{x} = \text{Dec}(\text{SES.ASA}[K])$. It then computes PIR request: $q \leftarrow \text{PIRQuery}(p, \mathbf{x})$ and forwards it to the cloud. The latter computes the response as $r \leftarrow \text{PIRResponse}(q)$. Finally the client decrypts r and the result is decrypted with $\text{SKE.Dec} = \text{pos}$, which is valid as long as $\text{pos} < n - m + 1$

Fig. 6: The Search and SrchToken algorithms of the S^3E protocol.

$24n + 4 * 2 * 2n = 40n$. The space cost of the solution based on suffix trees [9] can be further reduced to $20n$ by eliminating the need to store suffix links and parent pointers. In contrast, in S^3E we replace the space expensive suffix trees with suffix arrays and as such the index space cost is reduced from $20n$ bytes to $4n$ bytes.

For the sake of a clear comparison, for the storage space computation during the encryption of the index, be it suffix tree or suffix array, we exclude a per byte comparison and we assume a ciphertext comparison. The encryption of the index is based on the translation of the suffix tree to an encrypted dictionary. Thus, all the extra pointers of the suffix tree are excluded. Following the protocol from [9], the user encrypts $2n$ substrings which equal the number of edges of the suffix tree, n leaves and n characters of the original stream, resulting in $4n$ encryptions. In our solution thanks to the FM mapping the user does not need to send the original stream encrypted. It only sends the encrypted suffix array, plus two more n size arrays for the FM index construction; one for the F column of the index and one for the L column. In the end it uploads $4(n + \sum_{j=1}^k \max_c - f_{c_j})$ encrypted values to the cloud, in total.

Protocol	Search	Index [PS—CS]		Query [FR—LR]	VLS	Rounds	SL
CS[9]	$O(m+k)$	$20n$	$4n$	$\frac{m^2+\theta}{2\theta}$	$m+k$	✓	3 ✓+
FJKNRS[15]	$O(n)$	-	-	m	0	✓	1 ✓
FHV[16]	$O(n-m)$	-	-	m	0	✗	1 ✗
S ³ E	$O(m+k+p(n))$	$4n$	$4(n + \sum_{j=1}^k \max_c - f_{c_j})$	m	n	✓	2 ✗

Table 2: Comparison of existing substring searchable encryption protocols. Index space is further categorized in plaintext space index storage space (PS) and ciphertext space (CS). The overhead of [15] and [16] is undefined as the schemes do not take advantage of any auxiliary data structure for efficient substring search. For the query complexity we analyzed its size in terms of two separated phases: at the first round (FR) of the protocol and the last one (LR), in case of multiple rounds protocols. **VLS** denotes variable length substring search and **SL** defines whether or not the substring searchable protocol leaks the search pattern similarities. The + sign means that more than the search pattern is leaked.

For the query size we assume a block cipher of size θ and a substring query of size m . In [9] the substring is encrypted incrementally: for the substring 'abc' user encrypts separately $E(a)$, $E(ab)$, $E(abc)$. As such, for big substring queries as in DNA queries, the number of ciphertexts exceeds the number of the substring m . The total number of encryptions equals $1 + 2 + \dots + \frac{m}{\theta} = \frac{m^2+\theta}{2\theta}$ during the first round. At the last round the user asks for the positions of each character separately augmenting by a factor of m the substring size. In S³E the substring query has only per character encryptions of each character in the first round plus a PIR query of size n at the last round. Notice the PIR query can be computed in advance since it consists of multiple encryptions of 0 and encryption of 1. Our solutions also allows variable size substring queries, since the size of the substring query is decoupled from the scheme and can be defined online during the query phase as in [9].

For the rounds of communications S³E can return the substring search results in 2 rounds of communication; independent from the size of the stream or the substring query. During the first round the client sends an encrypted substring query and the cloud responds with encrypted addresses of the corresponding suffix array positions. At the second round the client issues a PIR query, the cloud replies with the PIR response and finally the client decrypts and learns the result. Our scheme with one less round of communication achieves better security than in [9] since we prevent the leakage of the search pattern and we achieve variable length substring queries with one more round of communication than in [16].

6 Performance

6.1 System

We implemented our substring searchable symmetric encryption scheme in order to see its practicality. For the client we used a machine running Ubuntu 14.04 with kernel version 3.19.0-29. The machine has 8MB RAM memory and is equipped with an INTEL Intel(R) Core(TM) i5-5200U CPU @ 2.20GHz processor with 4 cores. We instantiated the PRF, PRP and the public key encryption for the PIR query using Poly1305 + Salsa20, BLAKE2b [2] and NTRU cryptosystems, which outperform their competitors in computational efficiency. The most computational heavy operation on the client side is the computation of the index, which becomes the bottleneck of the total performance. In our benchmarks, since the code will vary in different machines, we choose to isolate the code in different modules in order to evaluate their relevant performance. Thus, the first module consists of operations that take place in cleartext data, the second module entails the cryptographic primitives used to encrypt the FM index and the third module consists also of all I/O operations in order to serialize the encrypted FM index. For our benchmarks we used two synthetic datasets of 24 and 4 alphabet size. We varied in either cases the size of the corpus in order to observe the feasibility results in these variations.

For the server implementation we parallelized the costly PIR computation by leveraging the computational power of GPUs. For our cluster we used 2x NVIDIA K1per K20X GPUs each with 6GB of memory, 2688 cores, 64 KB constant memory, 48 KB shared memory being able to achieve 250GB/sec memory throughput. The hardware setting enables a grid of $2147483647 \times 65535 \times 65535$ dimensions and maximum block dimension $1024 \times 1024 \times 64$ for thread parallelization. The warp size which denotes the maximum size of threads to coalesce memory accesses in hardware is 32. The multiple core processors of GPUs and the exploitation of the shared memory paradigm with tile programming, which carefully loads memory elements from global memory to the shared memory of each thread block, enable us to achieve fast matrix multiplications, which is the core server side bottleneck for substring queries responses.

6.2 Benchmarks

We present our benchmark results in this section starting with the cloud’s computational cost.

Cloud. We simulate the cloud with our GPU cluster setting and we measured the bottleneck of our substring symmetric searchable encryption protocol, which is the PIR response, which can be modeled as a vector (query) - matrix (data) multiplication. First we proceed with a naive matrix to vector multiplication implementation in a GPU environment. The results outperform by a factor of $\approx 4x$ non-parallelized single thread implementation. We further optimize the matrix multiplication operation by leveraging the common shared memory on GPU cores, in which threads of the same block on a grid have access to. The known as *tiled* technique tries to load in the shared memory tiles of data inputs that will be processed in parallel. That is, for a matrix multiplication we load to the shared memory rows and columns that will be processed for the computation of the final matrix. This careful memory treatment results in less memory latency for data transfer since accesses to the global and slow memory are substituted by fast shared memory accesses.

We present our results in tables 3, 4. From table 3 the GPU implementation, be it tiled or non-tiled, outperforms the single thread CPU implementation. We measured the PIR response time for different streams of data varying in $10^5, 5 \cdot 10^5, 10^6$ number of characters. In table 4 we benchmarked the time to copy memory from CPU to GPU and vice versa during the matrix multiplication for the PIR response. Thanks to the concurrent load of data blocks in tiles of the shared memory, multiple threads of the same block have access to the fast shared memory without the need of loading the data each time from the slow global memory. This results in a better computation to memory ratio, which dictates the number of float point operations with each access to the global memory.

Environment	#Characters		
	10^5	$5 \cdot 10^5$	10^6
Single Thread	0.804sec	3.993sec	8.547sec
GPU w/o Tiles	0.035sec	0.177sec	0.312sec
GPU with Tiles[32B]	0.024sec	0.120sec	0.300sec

Table 3: PIR response time

Environment	#Characters		
	10^5	$5 \cdot 10^5$	10^6
GPU w/o Tiles	0.199sec	0.521sec	0.987sec
GPU with Tiles[32B]	0.040sec	0.120sec	0.500sec

Table 4: GPU Memory Latency to copy data from and to GPU.

Client.

The computational cost of the client is dominated by the suffix array construction, its encryption and the query issue and decryption. We choose two different synthetic data sets of variable size; one with a vocabulary with 24 symbols and one with 4 symbols. In figure 5 we present our benchmark results, for computing the index¹, encrypting it and serialize it to the disk. The storage expansion of the index due to encryption is shown in table 5.

	10^5	$5 \cdot 10^5$	10^6
Encrypted Index	44MB	217MB	434MB
Plaintext Index	2.8MB	15MB	31MB

Table 5: Index Storage Capacity

The bottleneck at the client side is the PIR query since it has to be same size as the size of the data stream in order to hide the requested data cell. However this computational cost (refer to table 6)

¹ Code snippets can be found under: <http://pastebin.com/54HyZZBp>, <http://pastebin.com/MUpexJm0>

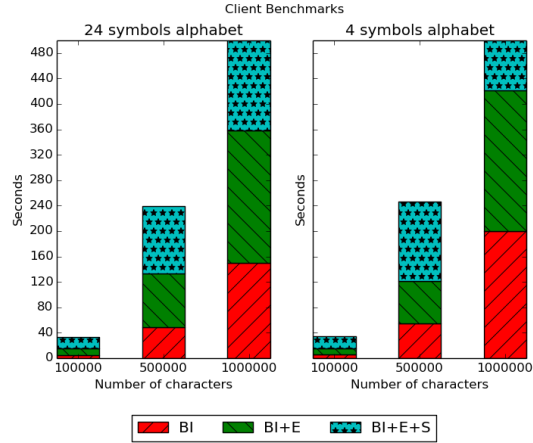


Fig. 7: BI:Build index, E:Encryption, S:Serialize

is independent of the query and can be computed offline. The user can compute encryption of 0s and encryption of 1 for the PIR query.

NTRU – EES401EP2	10^5	$5 \cdot 10^5$	10^6
	4.097s	10.368s	20.026s

Table 6: Query Complexity

7 Conclusion

We designed and analyzed a substring searchable symmetric encryption protocol S^3E , which achieves better storage performance, variable substring size and leaks less than state of the art work [9]. The idea of our protocol is to leverage the self-indexing mechanism of FM index, which stores only n integer positions of its substrings, without the need to store the mapping of substring in S . Our protocol is provably secure under the real-ideal indistinguishable simulation paradigm and its GPU parallelization at the cloud side demonstrates its real world practicality.

Bibliography

- [1] M. I. Abouelhoda, E. Ohlebusch, and S. Kurtz. Optimal exact string matching based on suffix arrays. In *In Proceedings of the Ninth International Symposium on String Processing and Information Retrieval*. Springer-Verlag, *Lecture Notes in Computer Science*, 2002.
- [2] J.-P. Aumasson, S. Neves, Z. Wilcox-O’Hearn, and C. Winnerlein. Blake2: Simpler, smaller, fast as md5. In *Proceedings of the 11th International Conference on Applied Cryptography and Network Security*, ACNS’13, pages 119–135, Berlin, Heidelberg, 2013. Springer-Verlag.
- [3] Bitglass. http://www.bitglass.com/company/news/press_releases/patentedencryption.
- [4] R. S. Boyer and J. S. Moore. A fast string searching algorithm. *Commun. ACM*, 20(10):762–772, Oct. 1977.
- [5] M. Burrows and D. J. Wheeler. A block-sorting lossless data compression algorithm. Technical report, 1994.
- [6] D. Cash, S. Jarecki, C. S. Jutla, H. Krawczyk, M. Rosu, and M. Steiner. Highly-scalable searchable symmetric encryption with support for boolean queries. In *Advances in Cryptology - CRYPTO 2013 - 33rd Annual Cryptology Conference, Santa Barbara, CA, USA, August 18-22, 2013. Proceedings, Part I*, pages 353–373, 2013.
- [7] Y.-C. Chang and M. Mitzenmacher. Privacy preserving keyword searches on remote encrypted data. In *Proceedings of the Third International Conference on Applied Cryptography and Network Security*, ACNS’05, pages 442–455, Berlin, Heidelberg, 2005. Springer-Verlag.
- [8] 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, Singapore, December 5-9, 2010. Proceedings*, pages 577–594, 2010.
- [9] M. Chase and E. Shen. Substring-searchable symmetric encryption. Cryptology ePrint Archive, Report 2014/638, 2014. <http://eprint.iacr.org/2014/638>.
- [10] Ciphercloud. <http://www.ciphercloud.com/technologies/encryption/>.
- [11] Ciphercloud. Q2-global-cloud-data-security-report. <http://pages.ciphercloud.com/rs/830-ILB-474/images/Q2-Global-Cloud-Data-Security-Report.pdf>.
- [12] R. Curtmola, J. Garay, S. Kamara, and R. Ostrovsky. Searchable symmetric encryption: Improved definitions and efficient constructions. In *Proceedings of the 13th ACM Conference on Computer and Communications Security*, CCS ’06, pages 79–88, New York, NY, USA, 2006. ACM.
- [13] E. De Cristofaro, S. Faber, and G. Tsudik. Secure genomic testing with size- and position-hiding private substring matching. In *Proceedings of the 12th ACM Workshop on Workshop on Privacy in the Electronic Society*, WPES ’13, pages 107–118, New York, NY, USA, 2013. ACM.
- [14] S. Devadas, M. van Dijk, C. W. Fletcher, L. Ren, E. Shi, and D. Wichs. Onion oram: A constant bandwidth blowup oblivious ram. Cryptology ePrint Archive, Report 2015/005, 2015. <http://eprint.iacr.org/2015/005>.
- [15] S. Faber, S. Jarecki, H. Krawczyk, Q. Nguyen, M. Rosu, and M. Steiner. Rich queries on encrypted data: Beyond exact matches. In *Computer Security - ESORICS 2015 - 20th European Symposium on Research in Computer Security, Vienna, Austria, September 21-25, 2015, Proceedings, Part II*, pages 123–145, 2015.
- [16] S. Faust, C. Hazay, and D. Venturi. Outsourced pattern matching. In F. Fomin, R. Freivalds, M. Kwiatkowska, and D. Peleg, editors, *Automata, Languages, and Programming*, volume 7966 of *Lecture Notes in Computer Science*, pages 545–556. Springer Berlin Heidelberg, 2013.
- [17] P. Ferragina and G. Manzini. Opportunistic data structures with applications. In *Proceedings of the 41st Annual Symposium on Foundations of Computer Science*, FOCS ’00, pages 390–, Washington, DC, USA, 2000. IEEE Computer Society.
- [18] R. Gennaro, C. Hazay, and J. S. Sorensen. Automata evaluation and text search protocols with simulation-based security. *J. Cryptology*, 29(2):243–282, 2016.
- [19] C. Gentry, K. A. Goldman, S. Halevi, C. S. Jutla, M. Raykova, and D. Wichs. Optimizing ORAM and using it efficiently for secure computation. In *Privacy Enhancing Technologies - 13th International Symposium, PETS 2013, Bloomington, IN, USA, July 10-12, 2013. Proceedings*, pages 1–18, 2013.
- [20] O. Goldreich, S. Goldwasser, and S. Micali. How to construct random functions. *J. ACM*, 33(4):792–807, Aug. 1986.

- [21] O. Goldreich and R. Ostrovsky. Software protection and simulation on oblivious rams. *J. ACM*, 43(3):431–473, May 1996.
- [22] G. H. Gonnet, R. A. Baeza-Yates, and T. Snider. Information retrieval. chapter New Indices for Text: PAT Trees and PAT Arrays, pages 66–82. Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1992.
- [23] C. Hazay and Y. Lindell. Efficient protocols for set intersection and pattern matching with security against malicious and covert adversaries. *J. Cryptology*, 23(3):422–456, 2010.
- [24] C. Hazay and T. Toft. Computationally secure pattern matching in the presence of malicious adversaries. In *Advances in Cryptology-ASIACRYPT 2010*, pages 195–212. Springer Berlin Heidelberg, 2010.
- [25] Hitachi. Searchable encryption: A technology for supporting secure application. http://www.hitachi.com/rd/portal/contents/story/searchable_encryption/index.html.
- [26] J. Kärkkäinen and P. Sanders. Simple linear work suffix array construction. In *Automata, Languages and Programming, 30th International Colloquium, ICALP 2003, Eindhoven, The Netherlands, June 30 - July 4, 2003. Proceedings*, pages 943–955, 2003.
- [27] R. M. Karp and M. O. Rabin. Efficient randomized pattern-matching algorithms. *IBM J. Res. Dev.*, 31(2):249–260, Mar. 1987.
- [28] D. E. Knuth, J. H. M. Jr., and V. R. Pratt. Fast pattern matching in strings. *SIAM J. Comput.*, 6(2):323–350, 1977.
- [29] Y. Lindell. How to simulate it - A tutorial on the simulation proof technique. *IACR Cryptology ePrint Archive*, 2016:46, 2016.
- [30] U. Manber and G. Myers. Suffix arrays: A new method for on-line string searches. In *Proceedings of the First Annual ACM-SIAM Symposium on Discrete Algorithms, SODA '90*, pages 319–327, Philadelphia, PA, USA, 1990. Society for Industrial and Applied Mathematics.
- [31] E. M. McCreight. A space-economical suffix tree construction algorithm. *J. ACM*, 23(2):262–272, Apr. 1976.
- [32] Mitsubishielectric. <http://www.mitsubishielectric.com/news/2013/0703.html>.
- [33] T. Moataz and E.-O. Blass. Oblivious substring search with updates. *Cryptology ePrint Archive, Report 2015/722*, 2015. <http://eprint.iacr.org/2015/722>.
- [34] T. Moataz, T. Mayberry, and E.-O. Blass. Constant communication oram with small blocksize. In *Proceedings of the 22Nd ACM SIGSAC Conference on Computer and Communications Security, CCS '15*, pages 862–873, New York, NY, USA, 2015. ACM.
- [35] D. Papadopoulos, C. Papamantou, R. Tamassia, and N. Triandopoulos. Practical authenticated pattern matching with optimal proof size. *Proc. VLDB Endow.*, 8(7):750–761, Feb. 2015.
- [36] V. Pappas, F. Krell, B. Vo, V. Kolesnikov, T. Malkin, S. G. Choi, W. George, A. Keromytis, and S. Bellovin. Blind seer: A scalable private dbms. In *Security and Privacy (SP), 2014 IEEE Symposium on*, pages 359–374. IEEE, 2014.
- [37] E. Shi, T.-H. H. Chan, E. Stefanov, and M. Li. Oblivious ram with $o((\log n)^3)$ worst-case cost. In *Advances in Cryptology-ASIACRYPT 2011*, pages 197–214. Springer Berlin Heidelberg, 2011.
- [38] Skyhighnetworks. <https://www.skyhighnetworks.com/cloud-encryption/>.
- [39] E. Stefanov and E. Shi. Oblivstore: High performance oblivious cloud storage. In *Security and Privacy (SP), 2013 IEEE Symposium on*, pages 253–267, May 2013.
- [40] E. Stefanov, E. Shi, and D. Song. Towards practical oblivious ram. *arXiv preprint arXiv:1106.3652*, 2011.
- [41] E. Ukkonen. On-line construction of suffix trees. *Algorithmica*, 14(3):249–260.
- [42] Virtualworks. <http://www.virtualworks.com/viaworks-enterprise-search/>.
- [43] P. Weiner. Linear pattern matching algorithms. In *Switching and Automata Theory, 1973. SWAT '08. IEEE Conference Record of 14th Annual Symposium on*, pages 1–11, Oct 1973.

A Cryptographic Primitives

A.1 Pseudorandom functions (PRF)

Let the family of all functions in the universe from a domain X to a range Y to be $Func[X, Y]$. A truly random function $f \xleftarrow{\$} Func[X, Y]$ is chosen randomly from the set of $Func$. The set of all these functions

is $|Y|^{|X|}$ (gigantic number). It is true that for any random function f with range size L chosen randomly from $\text{Func}[X, Y]$, $\Pr[f(x) = y] = 2^{-L}$. The randomness is not parametrized neither by the size of X and Y nor by the size of the domain. We define a pseudorandom function $f_k : X \rightarrow Y$ as a function from the set of all functions from X to Y as soon as a particular key k is fixed.

Definition 4. Let $\text{Func} = \{F : X \rightarrow Y\}$ be a function family for all functions F that map elements from the domain X to the range R . Then a PRF = $\{f_k : X' \rightarrow Y'\} \subseteq \text{Func}$ for $k \xleftarrow{\$} K$, where K is the key space.

The security of a PRF is modeled with a game which is known as *real or random* security game[20]. Intuitively, an adversary \mathcal{A} is given access to an oracle that on input x from a domain X , flips a coin $b \xleftarrow{\$} \{0, 1\}$ and if $b = 0$ then it outputs $y = f(x)$, for $f \in \text{Func}[X, Y]$, otherwise it outputs $y = f_k(x)$. \mathcal{A} issues queries to the oracle polynomially many times on input of the security parameter λ . Finally \mathcal{A} outputs a guess b' for the bit b .

The advantage of a probabilistic polynomial time algorithm \mathcal{A} in the PRF game is

$$\text{Adv}_{\mathcal{A}}^{\text{PRF}} = \Pr[b \xleftarrow{\$} \{0, 1\}; b' \leftarrow \mathcal{A}(y) : b' = b]$$

Definition 5. A PRF is computationally secure if all probabilistic polynomially time algorithms \mathcal{A} have advantage in the PRF game: $\frac{1}{2} + \epsilon(\lambda)$, for a negligible function ϵ on the security parameter λ .

A.2 Pseudorandom permutations (PRP)

A permutation is a bijective function where the domain and the range are equal. Similarly with the random functions, let $\text{Perm}[X]$ to be the set of all permutations for the domain X . Then a pseudorandom permutation (PRP) is a randomly chosen permutation from the set $\text{Perm}[X]$, keyed under a secret key k .

The advantage of a probabilistic polynomial time algorithm \mathcal{A} in the PRP game is

$$\text{Adv}_{\mathcal{A}}^{\text{PRP}} = \Pr[b \xleftarrow{\$} \{0, 1\}; b' \leftarrow \mathcal{A}(y) : b' = b]$$

Definition 6. A PRP is computationally secure if all probabilistic polynomially time algorithms \mathcal{A} have advantage in the PRP game: $\frac{1}{2} + \epsilon(\lambda)$, for a negligible function ϵ on the security parameter λ .

A.3 Symmetric Key Encryption

A symmetric key encryption scheme $\text{SKE} = \{\text{Gen}, \text{Enc}, \text{Dec}\}$ consists of three algorithms. Gen takes as input a security parameter λ and outputs the secret key sk . The probabilistic encryption algorithm Enc takes as input the secret key sk and a plaintext x from the plaintext space \mathcal{P} and outputs the ciphertext c . The decryption algorithm SKE.Dec takes as input a ciphertext from the ciphertext space \mathcal{C} and the secret decryption key sk and outputs the plaintext $x \in \mathcal{P}$. Correctness follows $\iff \forall \text{sk} \leftarrow \text{Gen}(1^\lambda), \text{SKE.Dec}(\text{Enc}(\text{sk}, x)) = x, \forall x \in \mathcal{P}$. Security is modeled with the standard game based indistinguishability experiment for polynomial probabilistic time adversary \mathcal{A} .

$\text{PrivK}_{\mathcal{A}, \text{SKE}}(\lambda)$:

- \mathcal{A} has access to the security parameter 1^λ .
- A key $\text{sk} \leftarrow \text{Gen}(1^\lambda)$ is generated and \mathcal{A} can learn encryptions of x of its choice $x \in \mathcal{S} \subset \mathcal{P}$.
- Eventually \mathcal{A} outputs x_0, x_1 where $|x_0| = |x_1|$. $b \xleftarrow{\$}$ and $\text{Enc}(x_b, \text{sk})$ is returned to \mathcal{A} .
- \mathcal{A} outputs its guess for b, b' .

If $b' = b$ \mathcal{A} succeeds and the experiment $\text{PrivK}_{\mathcal{A}, \text{SKE}}(\lambda) = 1$.

Definition 7. A symmetric encryption scheme SEK has indistinguishable encryptions if the probabilities $\Pr[\text{PrivK}_{\mathcal{A}, \text{SKE}}(\lambda) = 1] \leq \frac{1}{2} + \text{neg}(\lambda)$.

A.4 Private Information Retrieval

A single server, trapdoor, computational private information retrieval protocol consists of a one round communication interaction between a honest client, who holds a vector of identifiers \mathbf{x} and an untrusted cloud, who stores a $n \times l$ matrix \mathbf{M} . In the end the client fetches a desired position i from the vector \mathbf{x} and the server learns nothing. More specifically a CPIR consists of the following three algorithms:

- $\mathbf{q} \leftarrow \text{PIRQuery}(p, \mathbf{x})$: The user precomputes a vector of size n , which contains encryptions of 0 in all but the cell is interested in, which is an encryption of 1 with the public encryption key p .
- $r \leftarrow \text{PIRResponse}(\mathbf{q})$: The server computes the dot product of the n -length vector \mathbf{q} with the $n \times l$ matrix \mathbf{M} . Thanks to the underlying homomorphic encryption of the query, the dot product annihilates all rows but the asked one, without the server identifying the row, thanks to the indistinguishability property of the encryption algorithm.
- $a \leftarrow \text{PIRDecode}(s, r)$: The client decrypts the vector and recovers the item at the i^{th} row.

We define now the notion of security for a *PIR* scheme, which is derived from the semantic security of the underlying encryption scheme E .

Definition 8. A *PIR* scheme is computationally secure, if for any PPT adversary \mathcal{A} its view in two *PIR* queries $\mathbf{q}_x \leftarrow \text{PIRQuery}(p, \mathbf{x})$ and $\mathbf{q}_y \leftarrow \text{PIRQuery}(p, \mathbf{y})$ is computationally indistinguishable.

B Security Analysis

Game	Change	Indistinguishability Argument
Game ₀	Game ₀ = $\mathbf{Real}_{\mathcal{A}(\lambda)}^{\text{S}^3\text{E}}$	By definition
Game ₁	Replace F_{k_f}	Pseudorandomness of F_{k_f}
Game ₂	Replace Π_{k_π}	Pseudorandomness of Π_{k_π}
Game ₃	Replace $\text{SKE} = \{\text{Gen}, \text{Enc}, \text{Dec}\}$	Semantically secure $\text{SKE} = \{\text{Gen}, \text{Enc}, \text{Dec}\}$
Game ₄	Replace PIR	Semantically secure E
Game ₅	Game ₅ = $\mathbf{Ideal}_{\mathcal{A}, \mathcal{S}(\lambda)}^{\text{S}^3\text{E}}$	By definition

Table 7: Hybrid games

Theorem 2. Let $F_{k_f}, \Pi_{k_\pi}, \text{SKE} = \{\text{Gen}, \text{Enc}, \text{Dec}\}, \text{PIR}$ be a pseudorandom function, a pseudorandom permutation, a semantically secure symmetric encryption scheme and a secure *PIR* scheme respectively, then our substring searchable symmetric encryption scheme S^3E is adaptively \mathcal{L} -semantically secure.

Proof. In the $\mathbf{Real}_{\mathcal{A}(\lambda)}^{\text{S}^3\text{E}}$ game the adversary plays the role of the cloud and the challenger the role of the client. In the beginning the Challenger selects uniformly at random keys $\mathbf{k} = (k_f, r, k_\pi, k_e)$ for a PRF $F_{k_f, r} : \{0, 1\}^\lambda \times \{0, 1\}^\nu \rightarrow \{0, 1\}^\mu$, a PRP $\Pi_{k_\pi} : \{0, 1\}^\lambda \times \{0, 1\}^\nu \rightarrow \{0, 1\}^\nu$ and a symmetric encryption algorithm $\text{SKE} = \{\text{Gen}, \text{Enc}, \text{Dec}\}$. Upon receipt of a stream S of size n , the Challenger employs the $\text{SES} \leftarrow \text{PreProcess}(\mathbf{k}, S)$ as presented in figure 5 and forwards SES to \mathcal{A} . We distinguish between matching q^m and non-matching queries $q^{nm} : \mathbf{q} = \bigcup q^{nm} q^m$. Upon receiving the substring queries \mathbf{q} , the Challenger with F_{k_f} computes $\text{tk}_{T, S} = C[1], C[2], \dots, C[m] \leftarrow F_{k_f}(T[1..m]), F_r(T[m])$. We assume for the ease of readability that adversary issues only matching queries q^m . Finally \mathcal{A} receives $\mathbf{t} = (\text{tk}_1, \text{tk}_2, \text{tk}_3, \dots, \text{tk}_o)$ for each substring query.

Within a sequence of hybrid games we show the indistinguishable transformation of $\mathbf{Real}_{\mathcal{A}(\lambda)}^{\text{S}^3\text{E}}$ game to eventually the $\mathbf{Ideal}_{\mathcal{A}, \mathcal{S}(\lambda)}^{\text{S}^3\text{E}}$ game, which concludes the proof. The simulator \mathcal{S} computes the simulated encrypted index $\text{SES}^* = (\text{LLSet}^*, \text{LL}_c^*, \text{FM}^*)$ as follows:

- **Game₀**: This game is equivalent with the $\mathbf{Real}_{\mathcal{A}(\lambda)}^{\text{S}^3\text{E}}$ game.
- **Game₁**: This game behaves as the $\mathbf{Real}_{\mathcal{A}(\lambda)}^{\text{S}^3\text{E}}$ game with the difference that \mathcal{S} does not have access to S . The simulator through the \mathcal{L}_1 leakage function builds the substring encrypted structure SES as follows: We assume the existence of an algorithm $S \leftarrow \text{Build}(n', str)$, which takes as input $n' \in \mathbb{N}$ and the structure $str = \{c\}_{i=1}^{n'}, c \in \Sigma^*$ and outputs a bitstring of length n' , from a vocabulary Σ^* . Notice

that as in the real game the valid length of the original stream is not revealed and only the length of the string after the padding n' is leaked. \mathcal{S} selects uniformly at random keys $k = (k_f, r, k_\pi, k_e)$ for a PRF $F_{k_f} : \{0, 1\}^\lambda \times \{0, 1\}^\nu \rightarrow \{0, 1\}^\mu$, a PRP $\Pi_{k_\pi} : \{0, 1\}^\lambda \times \{0, 1\}^\nu \rightarrow \{0, 1\}^\nu$ and a symmetric encryption algorithm $SKE = \{\text{Gen}, \text{Enc}, \text{Dec}\}$ and runs $SES \leftarrow \text{PreProcess}(k, \text{Build}(\mathcal{L}_1))$. \mathcal{S} uses F_{k_f} to evaluate bit strings of length c_n $\mathcal{L}_2(q) = c_n, str$.

- **Game₂**: This game behaves similarly with **Game₁**, but we replace the F_{k_f} with a real random function which is evaluated through access to an oracle $\mathcal{O}^{\text{RF}}(\lambda, \mu, \nu)$.
- **Game₃**: This game behaves similarly with **Game₂**, but we replace the Π_{k_π} with a real random permutation which is evaluated through access to an oracle $\mathcal{O}^{\text{RP}}(\lambda, \nu)$.
- **Game₄**: In **Game₄** we replace the semantically secure $SKE = \{\text{Gen}, \text{Enc}, \text{Dec}\}$ with real random values by querying an oracle $\mathcal{O}^{\text{RE}}(\lambda)$.
- **Game₅**: In **Game₅** is identical with **Game₄** apart the behavior of the PIR in which we replace the underlying asymmetric homomorphic encryption algorithm E with real random values, through access to an oracle $\mathcal{O}^{\text{RHE}}(\lambda)$.

We write $\text{Game}_i \approx \text{Game}_j$ to denote that the view of probabilistic polynomial time adversary \mathcal{A} is indistinguishable between the output of **Game_i** and **Game_j**. $\text{Game}_0 = \mathbf{Real}_{\mathcal{A}(\lambda)}^{\text{S}^3\text{E}}$ by definition, $\text{Game}_1 \approx \text{Game}_0$ as long as no collisions happen to the evaluation of F_{k_f} , Π_{k_π} , $SKE = \{\text{Gen}, \text{Enc}, \text{Dec}\}$ or E , $\text{Game}_2 \approx \text{Game}_1$ as long as F_{k_f} is indistinguishable from real random function, $\text{Game}_3 \approx \text{Game}_2$ thanks to the indistinguishable output of Π_{k_π} from real random permutations, $\text{Game}_4 \approx \text{Game}_3$ because of the semantically secure $SKE = \{\text{Gen}, \text{Enc}, \text{Dec}\}$, $\text{Game}_5 \approx \text{Game}_4$ from the indistinguishable outputs of the symmetric encryption scheme E and finally $\text{Game}_5 = \mathbf{Ideal}_{\mathcal{A}, \mathcal{S}(\lambda)}^{\text{S}^3\text{E}}$ by definition.