

# Searchable Symmetric Encryption: Optimal Locality in Linear Space via Two-Dimensional Balanced Allocations

Gilad Asharov\*      Moni Naor<sup>†</sup>      Gil Segev<sup>‡</sup>      Ido Shahaf<sup>‡</sup>

## Abstract

Searchable symmetric encryption (SSE) enables a client to store a database on an untrusted server while supporting keyword search in a secure manner. Despite the rapidly increasing interest in SSE technology, experiments indicate that the performance of the known schemes scales badly to large databases. Somewhat surprisingly, this is not due to their usage of cryptographic tools, but rather due to their poor *locality* (where locality is defined as the number of *non-contiguous* memory locations the server accesses with each query). The only known schemes that do not suffer from poor locality suffer either from an impractical *space overhead* or from an impractical *read efficiency* (where read efficiency is defined as the ratio between the number of bits the server reads with each query and the actual size of the answer).

We construct the first SSE schemes that simultaneously enjoy optimal locality, optimal space overhead, and nearly-optimal read efficiency. Specifically, for a database of size  $N$ , under the modest assumption that no keyword appears in more than  $N^{1-1/\log \log N}$  documents, we construct a scheme with read efficiency  $\tilde{O}(\log \log N)$ . This essentially matches the lower bound of Cash and Tessaro (EUROCRYPT '14) showing that any SSE scheme must be sub-optimal in either its locality, its space overhead, or its read efficiency. In addition, even without making any assumptions on the structure of the database, we construct a scheme with read efficiency  $\tilde{O}(\log N)$ .

Our schemes are obtained via a *two-dimensional* generalization of the classic balanced allocations (“balls and bins”) problem that we put forward. We construct nearly-optimal two-dimensional balanced allocation schemes, and then combine their algorithmic structure with subtle cryptographic techniques.

---

A preliminary version of this work appeared in *Proceedings of the 48th Annual ACM Symposium on Theory of Computing (STOC)*, pages 1101–1114, 2016.

\*IBM Research. Email: [gsasharo@us.ibm.com](mailto:gsasharo@us.ibm.com). This work was completed while the author was a post-doctoral researcher at the Hebrew University’s School of Computer Science and Engineering.

<sup>†</sup>Incumbent of the Judith Kleeman Professorial Chair, Department of Computer Science and Applied Mathematics, Weizmann Institute of Science, Rehovot 76100, Israel. Email: [moni.naor@weizmann.ac.il](mailto:moni.naor@weizmann.ac.il). Research supported in part by the Israel Science Foundation and by the Israeli Centers of Research Excellence (I-CORE) Program (Center No. 4/11).

<sup>‡</sup>School of Computer Science and Engineering, Hebrew University of Jerusalem, Jerusalem 91904, Israel. Email: [{segev,ido\\_shahaf}@cs.huji.ac.il](mailto:{segev,ido_shahaf}@cs.huji.ac.il). Supported by the Israel Science Foundation (Grant No. 483/13), by the Israeli Centers of Research Excellence (I-CORE) Program (Center No. 4/11), and by the European Union’s Horizon 2020 Framework Program (H2020) via an ERC Grant (Grant No. 714253).

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Our Contributions . . . . .	2
1.2	Related Work . . . . .	2
1.3	Overview of Our Approach . . . . .	4
1.4	Paper Organization . . . . .	6
<b>2</b>	<b>Preliminaries</b>	<b>7</b>
2.1	Searchable Symmetric Encryption . . . . .	7
2.1.1	Functionality . . . . .	8
2.1.2	Locality and Read Efficiency . . . . .	8
2.1.3	Simulation-Based Security . . . . .	9
2.2	Additional Tools . . . . .	10
2.3	Negative Association . . . . .	11
2.4	Some Standard Tail Bounds . . . . .	12
<b>3</b>	<b>Two-Dimensional Balanced Allocations</b>	<b>13</b>
3.1	A One-Choice Allocation Scheme . . . . .	13
3.2	A Two-Choice Allocation Scheme . . . . .	15
3.2.1	Bounding the Maximal Load . . . . .	16
3.2.2	The Tightness of Our Analysis . . . . .	22
<b>4</b>	<b>SSE Schemes from Two-Dimensional Balanced Allocations</b>	<b>22</b>
4.1	Allocation Algorithms . . . . .	23
4.2	From Allocation Algorithms to SSE Schemes . . . . .	24
4.3	Extensions . . . . .	26
4.4	A Limitation of Allocation Algorithms . . . . .	27
<b>5</b>	<b>Improving the Cash-Tessaro Scheme</b>	<b>29</b>
	<b>References</b>	<b>32</b>

# 1 Introduction

Outsourcing data storage to remote servers is an extremely useful technology that has been adopted both by large organizations and by individual users. It offers great benefits, but at the same time, raises various concerns when dealing with sensitive data. In particular, in order to preserve the confidentiality of the data against an untrusted server, user-side symmetric encryption methods are typically employed prior to storing the data. As a result, operations as basic as keyword search become extremely expensive and sometimes even infeasible. This problem has motivated the cryptographic community to develop encryption methods that enable to search over symmetrically-encrypted data while not revealing sensitive information.

A searchable symmetric encryption (SSE) scheme [SWP00, CGK<sup>+</sup>06] is a mechanism that allows a client to store data on an untrusted server and later perform keyword searches: Given a keyword  $w$ , the client should be able to retrieve all documents that contain  $w$ . First, the client encrypts its database and uploads it to the server. The client can then repeatedly query the server with various keywords. Informally, the security requirement asks that the server does not learn any information about keywords for which the client did not issue any queries.

A very productive line of research (see, for example, [SWP00, Goh03, CM05, CGK<sup>+</sup>06, CK10, vLSD<sup>+</sup>10, CGK<sup>+</sup>11, KO12, KPR12, CJJ<sup>+</sup>13, KO13, KP13, CJJ<sup>+</sup>14, CT14, CGP<sup>+</sup>15] and the references therein) has been devoted to the construction of searchable symmetric encryption schemes. However, implementations and experiments with real-world databases [CJJ<sup>+</sup>13] indicate that the performance of the known schemes is quite disappointing and scales badly to large databases. Somewhat surprisingly, it turns out that the main bottleneck is in fact not the cryptographic processing of the data, but rather lower-level issues resulting from the inefficient memory layouts required by these schemes.

Specifically, the common drawback in the known schemes is poor *locality*: the server has to access a rather large number of *non-contiguous* memory locations with each query. The only known schemes that do not suffer from poor locality suffer either from an impractical *space overhead* (i.e., their encrypted databases are much larger than the original databases), or from an impractical *read efficiency* (i.e., they read much more data than needed for answering each query).<sup>1</sup>

This state of affairs naturally poses the challenge of constructing a searchable symmetric encryption scheme that simultaneously enjoys optimal space overhead, locality, and read efficiency. However, Cash and Tessaro [CT14] elegantly proved that obtaining such optimal efficiency guarantees is in fact impossible: any scheme must be sub-optimal in either its space overhead, its locality, or its read efficiency. Cash and Tessaro also presented a scheme that significantly improved the trade-off between the space overhead, locality, and read efficiency of the previously known schemes. However, although their scheme offers optimal read efficiency, it is far from optimal in both its space overhead and locality (see Table 1). This leads to the following fundamental problem, introducing a subtle combination of cryptographic and algorithmic challenges:

*Can we construct a searchable symmetric encryption scheme that simultaneously enjoys nearly-optimal space overhead, locality, and read efficiency?*

---

<sup>1</sup>In this work we rely on the notions of locality and read efficiency as formalized by Cash and Tessaro [CT14]. Specifically, the locality of a scheme is the maximal number of non-contiguous memory accesses that the server performs with each query, and the read efficiency of a scheme is the ratio between the number of bits the server reads with each query and the actual size of the answer. We refer the reader to Section 2.1 for the formal definitions.

## 1.1 Our Contributions

We construct searchable symmetric encryption schemes that simultaneously enjoy nearly-optimal space overhead, locality, and read efficiency. Our main results are obtained by basing searchable symmetric encryption on a *two-dimensional* generalization of the classic balanced allocations problem (also known as the “balls and bins” problem) that we put forward. We construct efficient two-dimensional balanced allocation schemes, and by combining them with subtle cryptographic techniques we obtain significant improvements over the previously known searchable symmetric encryption schemes.

Following the line of research on searchable symmetric encryption, we assume that a database is represented as a collection  $\text{DB} = \{\text{DB}(w_1), \dots, \text{DB}(w_{n_W})\}$ , where  $w_1, \dots, w_{n_W}$  are distinct keywords, and  $\text{DB}(w)$  is the list of all documents (or document identifiers) that contain the keyword  $w$ . We denote by  $N = \sum_{i=1}^{n_W} |\text{DB}(w_i)|$  the size of the database, and for each keyword  $w$  we denote by  $n_w = |\text{DB}(w)|$  the number of documents containing  $w$ . We consider the standard unit-cost word-RAM model, and assume for simplicity that keywords and document identifiers are represented using a constant number of machine words (thus we measure space usage in machine words)<sup>2</sup>.

We construct three searchable symmetric encryption schemes whose efficiency guarantees are summarized in Theorem 1.1 and in Table 1.<sup>3</sup> Whereas our first and second schemes are based on the above mentioned two-dimensional balanced allocations problem, our third scheme follows a different approach. For our third scheme we show that the recent scheme of Cash and Tessaro [CT14] can be modified to have optimal locality without hurting its space overhead or read efficiency.

**Theorem 1.1.** *Assuming the existence of one-way functions, there exist searchable symmetric encryption schemes offering the following guarantees for databases of size  $N$ :*

1. *Space  $O(N)$ , locality  $O(1)$ , and read efficiency  $\tilde{O}(\log N)$  without any assumptions on the structure of the database.*
2. *Space  $O(N)$ , locality  $O(1)$ , and read efficiency  $\tilde{O}(\log \log N)$  assuming that no keyword appears in more than  $N^{1 - \frac{1}{\log \log N}}$  documents.*
3. *Space  $O(N \log N)$ , locality  $O(1)$ , and read efficiency  $O(1)$  without any assumptions on the structure of the database.*

In Table 1 we compare the efficiency guarantees of our schemes to those of the previously known schemes that have less than a polynomial space overhead<sup>4</sup>. Compared to those schemes, our first and second schemes show that optimal space and locality can be achieved together with a nearly-optimal read efficiency. Our third scheme shows that by slightly increasing the space overhead, it is possible to achieve optimal locality and read efficiency.

## 1.2 Related Work

**Searchable symmetric encryption.** The notion of searchable symmetric encryption was put forward by Song, Wagner and Perrig [SWP00] who suggested several practical constructions. Formal notions of security and constructions satisfying them were later provided by Curtmola, Garay, Kamara, and Ostrovsky [CGK<sup>+</sup>06, CGK<sup>+</sup>11]. Additional work in this line of research developed

---

<sup>2</sup>The unit cost word-RAM model is considered the standard model for analyzing the efficiency of data structures (see, for example, [DP08, Hag98, HMP01, Mil99, PP08] and the references therein).

<sup>3</sup>Throughout this paper, for any function  $f(n)$  we let  $\tilde{O}(f(n))$  denote  $O(f(n)(\log f(n))^c)$  for some constant  $c$ .

<sup>4</sup>There are various schemes (e.g. [CGK<sup>+</sup>06, CK10, vLSD<sup>+</sup>10, KO12, KP13]) that have a *polynomial* space overhead, and offer constant locality and read efficiency. However, for concreteness we compare our schemes only to those with a potentially practical space overhead.

Scheme	Space	Locality	Read Efficiency
[CGK <sup>+</sup> 06, KPR12, CJJ <sup>+</sup> 13]	$O(N)$	$O(n_w)$	$O(1)$
[CT14]	$O(N \log N)$	$O(\log N)$	$O(1)$
<b>This work I</b>	$O(N)$	$O(1)$	$\tilde{O}(\log N)$
<b>This work II</b>	$O(N)$	$O(1)$	$\tilde{O}(\log \log N)$
<b>This work III</b>	$O(N \log N)$	$O(1)$	$O(1)$
Lower bound [CT14]	$\omega(N)$	$O(1)$	$O(1)$

**Table 1:** The efficiency guarantees of our schemes and of the previously known schemes that have less than a polynomial space overhead. Recall that we denote by  $N$  the size of the underlying database, and by  $n_w$  the number of documents that contain each keyword  $w$  (note that  $n_w$  may be as large as  $N$ ). Our second scheme is based on the modest assumption that no keyword appears in more than  $N^{1-1/\log \log N}$  documents.

searchable symmetric encryption schemes with various efficiency properties [Goh03, CM05, CK10, vLSD<sup>+</sup>10, CT14], support for data updates [KPR12, KP13, CJJ<sup>+</sup>14], authenticity [KO13], and support for more advanced searches [CJJ<sup>+</sup>13].

The known constructions can be roughly divided into two, somewhat orthogonal, underlying approaches. The first approach (see, for example, [CGK<sup>+</sup>06, KPR12, CJJ<sup>+</sup>13] and the references therein) provides schemes with linear space and constant read efficiency, but with very poor locality. Given a database  $\text{DB} = \{\text{DB}(w_1), \dots, \text{DB}(w_{n_w})\}$  of size  $N$ , the idea underlying this approach is to allocate an array of size roughly  $N$ , and to uniformly map the  $N$  elements of the database into the  $N$  entries of the array (one element per entry). For efficiently recovering a list  $\text{DB}(w)$  given a keyword  $w$  (i.e., recovering the list of documents that contain  $w$ ), each document identifier in the list  $\text{DB}(w)$  is stored in the array together with a pointer to the next document identifier in the list. Thus, since the elements are uniformly mapped into the array, recovering  $\text{DB}(w)$  requires the server to access essentially  $n_w = |\text{DB}(w)|$  random locations in the array.

The second approach (see, for example, [CGK<sup>+</sup>06, CK10, vLSD<sup>+</sup>10, CGK<sup>+</sup>11, KO13, KP13] and the references therein) provides schemes with optimal locality and read efficiency, but with a significant space overhead. Given a database  $\text{DB} = \{\text{DB}(w_1), \dots, \text{DB}(w_{n_w})\}$  of size  $N$ , the idea underlying this approach is to allocate a sufficiently large array, and to uniformly map each list  $\text{DB}(w)$  into a contiguous interval of length  $n_w = |\text{DB}(w)|$  in the array, without any overlaps between different lists. For efficiently recovering a list  $\text{DB}(w)$  given a keyword  $w$ , the server needs to access only a single random location, and then read  $n_w$  consecutive entries (thus resulting in optimal locality and read efficiency). However, since the locations of the lists in the array reveal information of the structure of the underlying database, some padding must be applied for hiding information on the lengths of the lists (e.g., padding each list according to the length of the longest list), thus resulting in a *polynomial* space overhead. Very recently, Cash and Tessaro showed that the space overhead can in fact be reduced from polynomial to logarithmic, at the cost of increasing the locality from constant to logarithmic (see Table 1).

As mentioned above, the work of Cash and Tessaro [CT14] explored the trade-off between locality, space overhead, and read efficiency. More specifically, instead of considering read efficiency directly, they considered the measure of *overlapping reads*, where a scheme has  $\alpha$ -overlapping reads if reads for a new search overlaps in at most  $\alpha$  bits with all previous reads. As noted by Cash and Tessaro, a (somewhat weak) read efficiency requirement is implicit in the condition on overlapping reads, and they proved a super-linear space lower bound that ties together locality and overlapping reads

(see [CT14] for further details).

**Balanced allocations.** Our two-dimensional balanced allocations problem is a novel generalization of the classical balanced allocations problem (the “balls and bins” problem), and may be of independent interest. The classical balanced allocations problem has seen a large number of generalizations over the years, and we refer the reader to [MU05, MRS00] for an overview of this problem and of some of its generalizations. A multi-dimensional generalization of this problem was first introduced by Broder and Mitzenmacher [BM05]. They considered that task of throwing balls into bins, where each ball is a random  $D$ -dimensional 0-1 vector of weight  $f$  (i.e., each vector has exactly  $f$  non-zero entries, and is chosen uniformly from all  $\binom{D}{f}$  such vectors). Broder and Mitzenmacher were interested in the average load in each dimension for each bin. The multi-dimensional generalization that we consider, as we discuss in Section 1.3, seems completely different.

### 1.3 Overview of Our Approach

The problem of search on encrypted data offers an interesting challenge in the interface between data structure design and cryptography. We provide an overview of the main ideas underlying our schemes and in order to emphasize the new ones, we focus here on the first and second schemes, as these are based on a completely new approach for constructing searchable symmetric encryption schemes. Our third scheme follows a more standard approach (based on that of Cash and Tessaro [CT14]), and the reader is referred to Section 5 for the main ideas underlying that scheme.

As discussed in Section 1.2, the known schemes can be roughly divided into two, somewhat orthogonal, underlying approaches. The first approach provides schemes with linear space and constant read efficiency, but with very poor locality. The second approach provides schemes with optimal locality and read efficiency, but with a significant space overhead. Our approach can be viewed as a subtle mixture of these two approaches that enables us to enjoy their advantages while mitigating their disadvantages.

Specifically, given a database  $DB = \{DB(w_1), \dots, DB(w_{n_w})\}$ , recall that the first approach completely ignores the structure of the lists, and places each document in a random location (leading to locality  $n_w = |DB(w)|$ ). The second approach preserves the structure of the lists and places the documents of each list in consecutive locations (leading to optimal locality and read efficiency, but with an impractical space overhead). Our approach preserves the structure of the lists, but does not place documents from the same list in consecutive locations. Instead, we make sure that documents from the same list are placed sufficiently near each other, yet in a random manner allowing documents from other lists to be placed between them. This flexibility enables us to enjoy constant locality in linear space, while only slightly affecting the read efficiency.

In what follows we begin by describing a *one-choice* scheme that follows this approach, and then describe a more efficient *two-choice* scheme. Throughout the exposition, we focus both on the algorithmic aspect of our schemes, introducing the *two-dimensional balanced allocations* problem, and on their cryptographic aspect.

**Warm-up: A one-choice scheme.** Given a database  $DB = \{DB(w_1), \dots, DB(w_{n_w})\}$  of size  $N$ , our one-choice scheme allocates an array of “bins”, and constructs an encrypted database as follows: For each keyword  $w$  we compute a hash value  $h(w)$ , and place the  $i$ th document from the list  $DB(w)$  in the bin  $h(w) + i - 1$ . That is, the first document containing  $w$  is placed in the bin  $h(w)$ , the second document containing  $w$  is placed in the bin  $h(w) + 1$ , and so on.

Once this process is completed for all keywords, note that each bin may contain more than a single document. Specifically, each bin may contain several documents corresponding to different

keywords. This can naturally be viewed as a two-dimensional generalization of the classic “balls and bins” problem, where one dimension consists of the keywords and the additional dimension consists of the structure of their lists (i.e., the lengths of their lists).

Equipped with this view of the process, we are interested in upper bounding its maximal load (that is, the maximal number of documents contained in each bin). However, this random process introduces new challenges when compared to its one-dimensional variant: The locations of the documents in the array are *not* independent. Specifically, given that a certain bin has many documents, it is rather likely that its adjacent bins also have many documents. Nevertheless, by carefully analyzing the dependencies that the lists introduce, we are able to bound its maximal load. Specifically, we show that for any database of size  $N$ , if we allocate an array of  $N/\tilde{O}(\log N)$  bins each of size  $\tilde{O}(\log N)$  (we make sure that the overall space is linear in  $N$ ), then with all but a negligible probability there are no overflowing bins.

This enables us to transform this algorithmic template into a searchable encryption scheme as follows. First, we pad each bin to contain exactly  $\tilde{O}(\log N)$  documents by adding “dummy” documents when needed, and we uniformly shuffle the documents in each bin. Then, the documents in each list  $\text{DB}(w)$  are each encrypted using an encryption key that is derived from the keyword  $w$  (say, by using a pseudorandom function that also generates the hash value  $h(w)$ ). We use an encryption scheme that produces *pseudorandom ciphertexts*, and has an *elusive and verifiable range* (see Section 2.2).

In terms of functionality, by providing the server with the value  $h(w)$  and with the decryption key corresponding to  $w$ , the server can recover the entire list  $\text{DB}(w)$  as follows. The server goes to the bin  $h(w)$ , and tries decrypting all of the ciphertexts that are stored there. Since we use an encryption scheme with an elusive and verifiable range, the server will be successful in decrypting exactly one of these ciphertexts, and this is the first document that contains  $w$ . The server then proceeds to the bin  $h(w) + 1$  and so on until it reaches a bin in which no ciphertext decrypts.

In terms of efficiency, the locality is optimal: The server goes to the bin  $h(w)$ , and from that point it only uses *contiguous* memory access. In addition, the read efficiency is  $\tilde{O}(\log N)$ , since for retrieving each document the server reads a bin containing  $\tilde{O}(\log N)$  encrypted documents. This yields a scheme with space  $O(N)$ , locality  $O(1)$ , and read efficiency  $\tilde{O}(\log N)$ .

Finally, the security of the scheme is based on the observation that the choices of bins for documents from different lists are statistically-close to being completely independent. Therefore, the access pattern that the server sees (i.e., the values  $h(w)$  for different keywords) does not reveal any unnecessary information on the database.

**An exponential improvement: A two-choice scheme.** In the classic balanced allocations problem it is well known that the “two-choice paradigm” [ABK<sup>+</sup>99, MU05] exponentially improves the maximal load for various ranges of the parameters<sup>5</sup>. This motivates us to consider a two-choice generalization of the above-described one-choice scheme.

Given a database  $\text{DB} = \{\text{DB}(w_1), \dots, \text{DB}(w_{n_w})\}$  of size  $N$ , our two-choice scheme allocates an array of bins and constructs an encrypted database as follows. For each keyword  $w$  we compute two independent hash values,  $h_1(w)$  and  $h_2(w)$ , and place all documents containing  $w$  in consecutive bins (as in our one-choice scheme) starting either from the bin  $h_1(w)$  or from the bin  $h_2(w)$  (but not from both). We choose between  $h_1(w)$  and  $h_2(w)$  as the starting point based on the maximal load of the bins in the relevant range (i.e., we choose the one with the smaller maximal load at the

---

<sup>5</sup>Roughly speaking, when throwing about  $N$  balls into  $N$  bins by placing each ball in the currently least loaded bin out of two possible choices, the maximal load is  $O(\log \log N)$  with high probability (as opposed to a maximal load of  $\Omega(\log N / \log \log N)$  in the one-choice process).

time of insertion)<sup>6</sup>.

This two-choice variant of our two-dimensional balanced allocations problem turns to be significantly more challenging to analyze (even when compared to the two-choice *one-dimensional* problem). Nevertheless, we are able to show that under the modest assumption that no keyword appears in more than  $N^{1-1/\log \log N}$  documents, we indeed obtain an exponential improvement: if we allocate an array of  $N/\tilde{O}(\log \log N)$  bins each of size  $\tilde{O}(\log \log N)$  (we again make sure that the overall space is linear in  $N$ ), then with all but a negligible probability there are no overflowing bins. Moreover, we show that the assumption that no keyword appears in more than  $N^{1-1/\log \log N}$  documents is in fact essential for obtaining such an exponential improvement.

Our analysis builds upon and generalizes the layered induction technique of Azar et al. [ABK<sup>+</sup>99]. Since our scheme places all elements of a list according to two random choices that are made for the *entire* list (and not for each element), this process introduces many dependencies between elements, as well as various dependencies between, say, the loads of any two consecutive bins. Such dependencies make the process significantly more challenging to analyze.

We once again transform this algorithmic template into a searchable encryption scheme. First, we pad each bin to contain exactly  $\tilde{O}(\log \log N)$  documents by adding “dummy” documents when needed, and we uniformly shuffle the documents in each bin. Then, the documents in each list  $\text{DB}(w)$  are each encrypted using an encryption key that is derived from the keyword  $w$  (by using a pseudorandom function as above). As in our one-choice scheme, we use an encryption scheme that produces *pseudorandom ciphertexts*, and has an *elusive and verifiable range*.

In terms of functionality, note that by providing the server with the values  $h_1(w)$  and  $h_2(w)$ , as well as with the number  $n_w$  of documents that contain the keyword  $w$ , the server can simply send back to the client the content of the  $n_w$  consecutive bins starting from the location  $h_1(w)$  and the content of the  $n_w$  consecutive bins starting from the location  $h_2(w)$ . The client can then decrypt the content of these bins and recover the list  $\text{DB}(w)$  using a decryption key that is derived from the keyword  $w$ . We note that we use an additional linear-space hash table for enabling the server to recover the value  $n_w$  on its own, and we emphasize that in this scheme we do not allow the server to decrypt the content of the bins on its own. This is due to the fact that whether the list  $\text{DB}(w)$  is stored starting from  $h_1(w)$  or  $h_2(w)$  may reveal unnecessary information on the structure of the database.

In terms of efficiency, as in our one-choice scheme we obtain optimal locality in linear space, but here we improve the read efficiency from  $\tilde{O}(\log N)$  to  $\tilde{O}(\log \log N)$ : for retrieving each document the server reads the content of two bins, each of which contains  $\tilde{O}(\log \log N)$  encrypted documents.

Finally, the security of the scheme is based on the observation that the two possible starting locations of documents from different lists in the database are statistically-close to being completely independent. Their actual starting locations are far from being independent, but this information is not revealed to the server since it cannot decrypt the content of the bins. Therefore, the access pattern that the server sees (i.e., the values  $h_1(w)$  and  $h_2(w)$  for various keywords  $w$ ) does not reveal any unnecessary information on the database.

## 1.4 Paper Organization

The remainder of this paper is organized as follows. In Section 2 we introduce the formal definition of symmetric searchable encryption schemes, as well as present various tools and tail bounds that are used in our constructions and proofs. Then, in Section 3 we put forward the two-dimensional balanced allocation problem, and then present and analyze a one-choice allocation scheme and a two-choice allocation scheme. In Section 4 we present a unified framework for constructing searchable

---

<sup>6</sup>In fact, our actual allocation rule is a bit more subtle, and we refer the reader to Section 3.2 for the specific details.



symmetric encryption schemes from two-dimensional allocation scheme, obtaining our first and second schemes as specific instantiations. Finally, in Section 5 we present our third scheme, improving the scheme of Cash and Tessaro [CT14].

## 2 Preliminaries

In this section we present the notation and basic definitions that are used in this work. We denote by  $\lambda \in \mathbb{N}$  the security parameter. For a distribution  $X$  we denote by  $x \leftarrow X$  the process of sampling a value  $x$  from the distribution  $X$ . Similarly, for a set  $\mathcal{X}$  we denote by  $x \leftarrow \mathcal{X}$  the process of sampling a value  $x$  from the uniform distribution over  $\mathcal{X}$ . For an integer  $n \in \mathbb{N}$  we denote by  $[n]$  the set  $\{1, \dots, n\}$ . A function  $\text{negl} : \mathbb{N} \rightarrow \mathbb{R}^+$  is **negligible** if for every constant  $c > 0$  there exists an integer  $N_c$  such that  $\text{negl}(n) < n^{-c}$  for all  $n > N_c$ . All logarithms in this paper are to the base of 2.

### 2.1 Searchable Symmetric Encryption

Let  $\lambda$  denote the security parameter. We let  $W = \{w_1, \dots, w_{n_W}\}$  denote the set of keywords, where  $n_W$  is polynomial in  $\lambda$ , and we assume that each keyword  $w_i$  can be represented using a constant number of machine words, each of length  $O(\lambda)$  bits, in the unit-cost RAM model. For each keyword  $w_i$ , we associate a list  $\text{DB}(w_i) = \{\text{id}_1, \dots, \text{id}_{n_i}\}$  of documents (or document identifiers) in which the keyword  $w_i$  appears. We assume that each identifier is of length  $O(\lambda)$  bits. We let  $\text{DB} = \{\text{DB}(w_1), \dots, \text{DB}(w_{n_W})\}$ , and let  $N = \sum_{i=1}^{n_W} n_i$  to denote the total number of keyword/identifier pairs in  $\text{DB}$ . Finally, let  $n_{\text{id}}$  denote the number of unique identifiers.

**Searchable symmetric encryption.** There are various different syntaxes for SSE schemes in the literature, where the main differences are the complexity of the interaction between the server and the client with each query. While some consider an interaction between the two that may consist of few rounds (where the server learns no information), others consider a single round interaction where the server also decrypts the data and sends the result to the client (and thus the server learn the output of each query). Usually, apart from the index of keywords/identifiers, the client uploads a sequence of encrypted documents, and the index is used to help the client locate the documents corresponding to its keyword.

Some of our schemes consist of two rounds of interaction, whereas some consist of a single round (and can thus be easily adapted to two rounds). In a one-round scheme, the client sends a token, the server performs some computation during which it learns the set of identifiers. It then fetches the corresponding (encrypted) documents from the database and sends them back to the client, which decrypts them and learns the result. In a two-round scheme, the server does not decrypt the set of identifiers on its own, and it sends the client some message, which the client resolves and learns the set of identifiers and send them back to the server. As in a one-round scheme, given the identifiers the server fetches the encrypted documents from the database and sends them back to the client.

In fact, two-rounds schemes can be viewed as having one round of interaction, if the server stores a database of keywords/documents instead of a database of keywords/identifiers. We note that the lower bound of [CT14] holds in both setting, since their lower bound relates only to the *structure* of the database. We proceed to the formal definitions. Although the following definition seems to relate to one-round schemes, it actually captures both options. This is because the second round is just the transmission of the identifiers, and fetching them from storage.

### 2.1.1 Functionality

A searchable symmetric encryption (SSE) scheme  $\Pi$  consists of algorithms (KeyGen, EDBSetup, TokGen, Search, Resolve) such that:

- $K \leftarrow \text{KeyGen}(1^\lambda)$ . The key generation algorithm KeyGen takes as input the security parameter  $1^\lambda$  and outputs a secret key  $K$ .
- $\text{EDB} \leftarrow \text{EDBSetup}(K, \text{DB})$ . The database-setup EDBSetup algorithm takes as input a key  $K$  and a database DB and outputs an encrypted database EDB.
- $(\tau, \rho) = \text{TokGen}(K, w)$ . The token generator algorithm is a deterministic algorithm that takes as input the private key  $K$  and the keyword  $w$  and outputs a token  $\tau$  to be sent to the server, and some internal state  $\rho$  for the resolving algorithm.
- $M = \text{Search}(\text{EDB}, \tau)$ . The searching algorithm Search takes as input the token  $\tau$  and the encrypted database EDB, and outputs some list  $M$  of results.
- $S = \text{Resolve}(\rho, M)$ . The resolving algorithm receives the list  $M$  and the state, and outputs the set of decrypted results  $L$ .

An SSE scheme for databases of size  $N = N(\lambda)$  is correct if for every PPT adversary  $\mathcal{A}$  there exists a negligible function  $\text{negl}(\cdot)$  such that the output of the experiment  $\text{correct}_{\mathcal{A}, \text{DB}}(1^\lambda)$  is 1 with probability  $1 - \text{negl}(\lambda)$  for all sufficiently large  $\lambda \in \mathbb{N}$ , where the experiment is defined as follows:

1. The adversary  $\mathcal{A}$  on input  $1^\lambda$  outputs a database DB of size  $N = N(\lambda)$  together with some state information **state**.
2. A key  $K \leftarrow \text{KeyGen}(1^\lambda)$  is chosen and the database is encrypted by computing  $\text{EDB} \leftarrow \text{EDBSetup}(K, \text{DB})$ .
3. The adversary  $\mathcal{A}$  is invoked on input  $(\text{EDB}, \text{state})$ , and can repeatedly issue queries  $w_i$ , where each query is answered as follows:
  - (a)  $(\tau_i, \rho_i) \leftarrow \text{TokGen}(K, w_i)$ ,  $M_i \leftarrow \text{Search}(\tau_i, \text{EDB})$  and  $S_i = \text{Resolve}(\rho_i, M_i)$ .
4. The output of the experiment is 1 if and only if for every query  $w_i$  it holds that  $S_i = \text{DB}(w_i)$ .

We note that in case where the database is an index of keywords/documents, then the above syntax is a one-round protocol. In case where the database is an index of keyword/identifiers, then this syntax corresponds to a two-round protocol, where the algorithms (TokGen, Search, Resolve) describes the first round, and in which in the second round has the following fixed structure: The client sends  $S$ , gets back the corresponding documents and then decrypts them. The single-round syntax for keyword/identifiers index is a sub-case of the above, where there is no Resolve algorithm, no secret state  $\rho$ , and the algorithm Search simply outputs  $S$  (and this is known to the server). Recall that the server then fetches the documents from storage according to the identifiers, and the client decrypts the actual document.

### 2.1.2 Locality and Read Efficiency

Our notions of locality and read efficiency follow those introduced by Cash and Tessaro [CT14].

**Locality.** The search procedure of any SSE scheme can be decomposed into a sequence of contiguous reads from the encrypted database EDB. Specifically, assume that the Search algorithm does not get the encrypted database as input, but rather only gets oracle access to it. Each query to this oracle consists of some interval  $[a_i, b_i]$ , and the oracle replies with the words that are stored on those intervals of EDB. The Search algorithm is invoked on some token  $\tau$  and queries its oracle on some

interval  $[a_1, b_1]$ . It continues to compute the next intervals to read based on  $\tau$  and all previously read intervals from EDB. We denote these intervals by  $\text{ReadPat}(\text{EDB}, \tau)$ .

**Definition 2.1** (Locality). *An SSE scheme  $\Pi$  is  $d$ -local (or has locality  $d$ ) if for every  $\lambda$ ,  $\text{DB}$  and  $w \in \mathcal{W}$ ,  $K \leftarrow \text{KeyGen}(1^\lambda)$ ,  $\text{EDB} \leftarrow \text{EDBSetup}(K, \text{DB})$  and  $\tau \leftarrow \text{TokGen}(K, w)$  we have that  $\text{ReadPat}(\text{EDB}, \tau)$  consists of at most  $d$  intervals with probability 1.*

**Read efficiency.** The notion of locality alone is lacking, since the Search algorithm can read the whole database with a single interval. The notion of read efficiency measures the overall size of portion read by a search operation. For a given  $\text{DB}$  and  $w$ , we let  $\|\text{DB}(w)\|$  denote the number of words in the encoding of  $\text{DB}(w)$ .

**Definition 2.2** (Read efficiency). *An SSE scheme  $\Pi$  is  $r$ -read efficient (or has read efficiency  $r$ ) if for any  $\lambda$ ,  $\text{DB}$ , and  $w \in \mathcal{W}$ , we have that  $\text{ReadPat}(\tau, \text{EDB})$  consists of intervals of total length at most  $r \cdot \|\text{DB}(w)\|$  words.*

### 2.1.3 Simulation-Based Security

The standard security definition for SSE schemes follows the ideal/real simulation paradigm. We consider both static and adaptive security, where the difference is whether the adversary chooses its queries statically (i.e., before seeing any token), or in an adaptive manner (i.e., the next query may be a function of the previous tokens). In both cases, some information is leaked to the server, which is formalized by letting the simulator receive the evaluation of some “leakage function” on the database itself and the real tokens. We start with the static case.

**The real execution.** The real execution is parameterized by the scheme  $\Pi$ , the adversary  $\mathcal{A}$ , and the security parameter  $\lambda$ . In the real execution the adversary is invoked on  $1^\lambda$ , and outputs a database  $\text{DB}$  and a list of queries  $\mathbf{w} = \{w_i\}_i$ . Then, the experiment invokes the key-generation algorithm and the database setup algorithms,  $K \leftarrow \text{KeyGen}(1^\lambda)$  and  $\text{EDB} \leftarrow \text{EDBSetup}(K, \text{DB})$ . Then, for each query  $w_i$  that the adversary has outputted, the token generator algorithm is run to obtain  $\tau_i = \text{TokGen}(w_i)$ . The adversary is given the encrypted database  $\text{EDB}$  and the resulting tokens  $\tau = \{\tau_i\}_{w_i \in \mathbf{w}}$ , and outputs a bit  $b$ .

**The ideal execution.** The ideal execution is parameterized by the scheme  $\Pi$ , a leakage function  $\mathcal{L}$ , the adversary  $\mathcal{A}$ , a simulator  $\mathcal{S}$  and the security parameter  $\lambda$ . In this execution, the adversary  $\mathcal{A}$  is invoked on  $1^\lambda$ , and outputs  $(\text{DB}, \mathbf{w})$  similarly to the real execution. However, this time the simulator  $\mathcal{S}$  is given the evaluation of the leakage function on  $(\text{DB}, \mathbf{w})$  and should output  $\text{EDB}, \tau$  (i.e.,  $(\text{EDB}, \tau) \leftarrow \mathcal{S}(\mathcal{L}(\text{DB}, \mathbf{w}))$ ). The execution follows by giving  $(\text{EDB}, \tau)$  to the adversary  $\mathcal{A}$ , which outputs a bit  $b$ .

Let  $\text{SSE-REAL}_{\Pi, \mathcal{A}}(\lambda)$  denote the output of the real execution, and let  $\text{SSE-IDEAL}_{\Pi, \mathcal{L}, \mathcal{A}, \mathcal{S}}(\lambda)$  denote the output of the ideal execution, with the adversary  $\mathcal{A}$ , simulator  $\mathcal{S}$  and leakage function  $\mathcal{L}$ . We now ready to define security of SSE:

**Definition 2.3** (static  $\mathcal{L}$ -secure SSE). *Let  $\Pi = (\text{KeyGen}, \text{EDBSetup}, \text{TokGen}, \text{Search})$  be an SSE scheme and let  $\mathcal{L}$  be a leakage function. We say that the scheme  $\Pi$  is static  $\mathcal{L}$ -secure searchable encryption if for every PPT adversary  $\mathcal{A}$ , there exists a PPT simulator  $\mathcal{S}$  and a negligible function  $\text{negl}(\cdot)$  such that*

$$|\Pr[\text{SSE-REAL}_{\Pi, \mathcal{A}}(\lambda) = 1] - \Pr[\text{SSE-IDEAL}_{\Pi, \mathcal{L}, \mathcal{A}, \mathcal{S}}(\lambda) = 1]| < \text{negl}(\lambda)$$

**Adaptive setting.** In the adaptive setting, the adversary is not restricted to specifying all of its queries  $\mathbf{w}$  in advance, but can instead choose its queries during the execution in an adaptive manner, depending on the encrypted database  $\text{EDB}$  and on the tokens that it sees. Let  $\text{SSE-REAL}_{\Pi, \mathcal{A}}^{\text{adapt}}(\lambda)$  denote the output of the real execution in this adaptive setting. In the ideal execution, the simulator  $\mathcal{S}$  is now an interactive Turing machine, which interacts with the experiment by responding to queries. First, the simulator  $\mathcal{S}$  is initially invoked on  $\mathcal{L}(\text{DB})$  and outputs  $\text{EDB}$ . Then, for every query  $w_i$  that  $\mathcal{A}$  may output, the function  $\mathcal{L}$  is invoked on  $\text{DB}$  and all previously queries  $\{w_j\}_{j < i}$  and the new query  $w_i$ , outputs some new leakage  $\ell(w_i)$  which is given to the simulator  $\mathcal{S}$ . The latter outputs some  $t_i$ , which is given back to  $\mathcal{A}$ , who may then issue a new query. At the end of the execution,  $\mathcal{A}$  outputs a bit  $b$ . Let  $\text{SSE-IDEAL}_{\Pi, \mathcal{L}, \mathcal{A}, \mathcal{S}}^{\text{adapt}}(\lambda)$  be the output of the ideal execution. The adaptive security of SSE is defined as follows:

**Definition 2.4** (adaptive  $\mathcal{L}$ -secure SSE). *Let  $\Pi = (\text{KeyGen}, \text{EDBSetup}, \text{TokGen}, \text{Search})$  be an SSE scheme and let  $\mathcal{L}$  be a leakage function. We say that the scheme  $\Pi$  is adaptive  $\mathcal{L}$ -secure searchable encryption if for every PPT adversary  $\mathcal{A}$ , there exists a PPT simulator  $\mathcal{S}$  and a negligible function  $\text{negl}(\cdot)$  such that*

$$\left| \Pr \left[ \text{SSE-REAL}_{\Pi, \mathcal{A}}^{\text{adapt}}(\lambda) = 1 \right] - \Pr \left[ \text{SSE-IDEAL}_{\Pi, \mathcal{L}, \mathcal{A}, \mathcal{S}}^{\text{adapt}}(\lambda) = 1 \right] \right| < \text{negl}(\lambda)$$

**The leakage function.** Following the standard notions of security for SSE we consider the leakage function  $\mathcal{L}_{\min}$  for one-round protocols and the leakage function  $\mathcal{L}_{\text{sizes}}$  for two-round protocols, where

$$\begin{aligned} \mathcal{L}_{\min}(\text{DB}, \mathbf{w}) &= (N, \{\text{DB}(w)\}_{w \in \mathbf{w}}), \\ \mathcal{L}_{\text{sizes}}(\text{DB}, \mathbf{w}) &= (N, \{|\text{DB}(w)|\}_{w \in \mathbf{w}}), \end{aligned}$$

and  $N = \sum_{w \in W} |\text{DB}(w)|$  is the size of the database. That is, both functions return the size of the database, and the difference between them is that the function  $\mathcal{L}_{\min}$  returns the actual documents that contain each keyword  $w \in \mathbf{w}$  that the adversary has queried, whereas the function  $\mathcal{L}_{\text{sizes}}$  returns only the number of such documents.

The leakage functions in the adaptive setting are defined analogously. That is, for a database  $\text{DB}$ , a set of “previous” queries  $\{w_j\}_{j < i}$ , and a new query  $w_i$ , we define

$$\begin{aligned} \mathcal{L}_{\min}^{\text{adap}}(\text{DB}, \{w_j\}_{j < i}, w_i) &= \begin{cases} N & \text{if } (\{w_j\}_{j < i}, w_i) = (\perp, \perp) \\ \text{DB}(w_i) & \text{otherwise} \end{cases} \\ \mathcal{L}_{\text{size}}^{\text{adap}}(\text{DB}, \{w_j\}_{j < i}, w_i) &= \begin{cases} N & \text{if } (\{w_j\}_{j < i}, w_i) = (\perp, \perp) \\ |\text{DB}(w_i)| & \text{otherwise} \end{cases}. \end{aligned}$$

## 2.2 Additional Tools

**Pseudorandom encryption schemes with elusive and verifiable range.** In our scheme, we will encrypt the identifiers of list with a different key, and we need to ensure that a decryption of some ciphertext is valid only when decrypting with the “correct” key. In addition, the client pads its lists using some random elements, and we want to ensure that by doing so it does not accidentally introduce valid ciphertexts for some keys, and that ciphertexts seems independent of their keys. We therefore use an encryption scheme with the following properties.

**Definition 2.5.** *Let  $(\text{Enc}, \text{Dec})$  be a private-key encryption scheme and denote the range of a key  $K \leftarrow \{0, 1\}^\lambda$  in the scheme by  $\text{Range}(K) \stackrel{\text{def}}{=} \{\text{Enc}_K(x)\}_{x \in \{0, 1\}^n}$ . Then,*

1. We say that  $(\text{Enc}, \text{Dec})$  has an elusive range if for every PPT machine  $\mathcal{A}$  there exists a negligible function  $\text{negl}(\cdot)$  such that

$$\Pr_{K \leftarrow \{0,1\}^\lambda} \left[ \mathcal{A}(1^\lambda) \in \text{Range}(K) \right] < \text{negl}(\lambda)$$

2. We say that  $(\text{Enc}, \text{Dec})$  has an efficiently verifiable range if there exists a PPT machine  $M$  such that  $M(1^\lambda, K, c) = 1$  if and only if  $c \in \text{Range}(K)$ .  
By convention, for every  $c \notin \text{Range}(K)$ , we have that  $\text{Dec}_K(c) = \perp$ .
3. We say that  $(\text{Enc}, \text{Dec})$  has pseudorandom ciphertexts if for every PPT adversary  $\mathcal{A}$  there exists a negligible function  $\text{negl}(\cdot)$  such that:

$$\left| \Pr \left[ \mathcal{A}^{\text{Enc}_K(\cdot)}(1^\lambda) = 1 \right] - \Pr \left[ \mathcal{A}^{\mathcal{R}(\cdot)}(1^\lambda) = 1 \right] \right| \leq \text{negl}(\lambda)$$

where  $\mathcal{R}$  is a probabilistic oracle that given any input outputs a freshly-sampled uniform value of the appropriate length (i.e., as the output length of  $\text{Enc}_K(\cdot)$ ), the probability on the left is taken over the choice  $K \leftarrow \{0,1\}^\lambda$  and the internal randomness of the algorithm  $\text{Enc}$ , and the probability on the right is taken over the randomness of the oracle  $\mathcal{R}$ .

We omit the  $\text{Gen}$  algorithm since the scheme is a symmetric-key one and thus we can assume w.l.o.g. that the key is a uniform string (i.e., the randomness for  $\text{Gen}$ ). We note that an encryption scheme that satisfy all these properties can be constructed from any pseudorandom function (and therefore from one-way function)  $\text{PRF}_K : \{0,1\}^m \rightarrow \{0,1\}^{m+\lambda}$  for  $K \in \{0,1\}^\lambda$  as  $\text{Enc}_K(x) = \langle r, \text{PRF}_K(r) \oplus x0^\lambda \rangle$ , where  $x \in \{0,1\}^m$ ,  $r \leftarrow \{0,1\}^\lambda$ , and  $x0^\lambda$  denotes the concatenation of  $x$  and  $0^\lambda$ . This scheme satisfy the first two notions [LP09], and it is easy to see that it has also pseudorandom ciphertexts.

**Static hash tables.** In our schemes we rely on static hash tables (also known as static dictionaries). These are data structures that given a set  $S$  can support lookup operations in constant time in the standard unit-cost word-RAM model. Specifically, a static hash table consists of a pair of algorithms denoted  $(\text{HTSetup}, \text{HTLookup})$ . The algorithm  $\text{HTSetup}$  gets as input a set  $S = \{(\ell_i, d_i)\}_{i=1}^k$  of pairs  $(\ell_i, d_i)$  of strings, where  $\ell_i \in \{0,1\}^s$  is the label and  $d_i \in \{0,1\}^r$  is the data. The output of this algorithm is a hash table  $\text{HT}(S)$ . The lookup algorithm  $\text{HTLookup}$  on input  $(\text{HT}(S), \ell)$  returns  $d$  if  $(\ell, d) \in S$ , and  $\perp$  otherwise.

There exist many constructions of static hash tables that use linear space (i.e.,  $O(k(r+s))$  bits) and answer lookup queries by reading a constant number of contiguous  $s$ -bit blocks and  $r$ -bit blocks (see, for example, [PR04, ANS10], and the many references therein).

## 2.3 Negative Association

Our proofs in this paper use the notion of negative association. Intuitively, random variables  $X_1, \dots, X_n$  are negatively associated if for every pair of disjoint sets  $S, T \subseteq [n]$ , if the variables  $(X_i)_{i \in S}$  are “high” then the variables  $(X_i)_{i \in T}$  are likely to be “low”. For a survey of negative association we refer the reader to the notes of Wajc [Waj17].

**Definition 2.6** (Negative Association [KS81, JDP83]). *Let  $X_1, \dots, X_n$  be random variables taking values in  $\mathbb{R}$ . We say that  $X_1, \dots, X_n$  are negatively associated if for every pair of disjoint sets  $S, T \subseteq [n]$ , and for every pair of monotonically-increasing function  $f : \mathbb{R}^{|S|} \rightarrow \mathbb{R}$  and  $g : \mathbb{R}^{|T|} \rightarrow \mathbb{R}$ , it holds that*

$$\mathbb{E} [f((X_i)_{i \in S}) \cdot g((X_j)_{j \in T})] \leq \mathbb{E} [f((X_i)_{i \in S})] \cdot \mathbb{E} [g((X_j)_{j \in T})] .$$

We make use of the following lemmas.

**Lemma 2.7** (Zero-one principle [DR98]). *If  $X_1, \dots, X_n$  are random variables, each taking values in  $\{0, 1\}$ , such that  $\sum_{i=1}^n X_i = 1$  with probability 1, then  $X_1, \dots, X_n$  are negatively associated.*

**Lemma 2.8** (Concatenation lemma [DR98]). *If  $X_1, \dots, X_n$  and  $Y_1, \dots, Y_m$  are negatively associated, and the vectors  $X = (X_1, \dots, X_n)$  and  $Y = (Y_1, \dots, Y_m)$  are independent, then  $X_1, \dots, X_n, Y_1, \dots, Y_m$  are negatively-associated.*

**Lemma 2.9** ([KS81, JDP83]). *Let  $X_1, \dots, X_n$  be negatively-associated random variables, let  $V_1, \dots, V_k \subseteq [n]$  be pairwise disjoint sets, and for each  $i \in [t]$  let  $f_i : \mathbb{R}^{|V_i|} \rightarrow \mathbb{R}$  be a monotone function. For a set  $S \subseteq [n]$ , denote  $X_S = (X_i)_{i \in S}$ . Then, the variables  $f_1(X_{V_1}), \dots, f_k(X_{V_k})$  are negatively associated.*

## 2.4 Some Standard Tail Bounds

Our proofs in this paper rely on the following standard tail-bound inequalities for sequences of random variables (see, for example, [MU05, DP09]).

**Lemma 2.10** (Multiplicative Chernoff's bound). *Suppose  $X_1, \dots, X_n$  are independent random variables taking values in  $\{0, 1\}$ . Let  $X$  denote their sum and let  $\mu = E[X]$  denote the sum's expected value. Then for any  $\delta > 0$ ,*

$$\Pr[X > (1 + \delta)\mu] \leq \left( \frac{e^\delta}{(1 + \delta)^{(1 + \delta)}} \right)^\mu.$$

**Lemma 2.11** (Hoeffding's bound for negatively-associated variables [DR98]). *Suppose  $X_1, \dots, X_n$  are negatively-associated random variables, where each  $X_i$  takes values in the interval  $[a_i, b_i]$ . Let  $X$  denote their sum and let  $\mu = E[X]$  denote the sum's expected value. Then for any  $t \geq 0$ ,*

$$\Pr[X - \mu \geq t] \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$

**Lemma 2.12.** *Suppose  $X_1, \dots, X_n$  and  $Y_1, \dots, Y_n$  are random variables such that each  $Y_i$  takes values in  $\{0, 1\}$ ,  $Y_i = Y_i(X_1, \dots, X_i)$  and  $\Pr[Y_i = 1 | X_1, \dots, X_{i-1}] \leq p$  for some  $p \in [0, 1]$ . Let  $Z_1, \dots, Z_n$  be independent Bernoulli trials with success probability  $p$ . Then for any non-negative coefficients  $c_1, \dots, c_n$  and a bound  $k$ ,*

$$\Pr\left[\sum_{i=1}^n c_i Y_i > k\right] \leq \Pr\left[\sum_{i=1}^n c_i Z_i > k\right].$$

Lemma 2.12 can be proved using a standard coupling argument (i.e., defining a joint probability distribution where always  $Y_i \leq Z_i$ ).

**Claim 2.13** (Bernstein's bound). *Let  $W_1, \dots, W_n$  be independent zero-mean random variables. Suppose that  $|W_j| \leq M$  for all  $j$ . Then, for any  $t > 0$ ,*

$$\Pr\left[\sum_{i=1}^n W_i > t\right] \leq \exp\left(-\frac{\frac{1}{2}t^2}{\sum_j E[W_j^2] + \frac{1}{3}Mt}\right).$$

**Corollary 2.14.** *Given Lemma 2.12's conditions, denote  $N = \sum_{i=1}^n c_i$  and  $M = \max c_i$ . Then,*

$$\Pr \left[ \sum_{i=1}^n c_i Y_i > 2Np \right] \leq \exp \left( -\frac{3}{8} \frac{N}{M} p \right).$$

**Proof.** Let  $Z_1, \dots, Z_n$  be the i.i.d. variables mentioned in the statement of Lemma 2.12. Set  $W_i = c_i(Z_i - p)$ . Then  $E[W_i] = 0$ ,  $|W_i| \leq c_i \max(1-p, p) \leq c_i$  and  $E[W_i^2] = p(1-p)c_i^2$ . Applying Lemma 2.12 and Bernstein's bound, we get

$$\begin{aligned} \Pr \left[ \sum_{i=1}^n c_i Y_i > 2Np \right] &\leq \Pr \left[ \sum_{i=1}^n c_i Z_i > 2Np \right] \\ &= \Pr \left[ \sum_{i=1}^n W_i > Np \right] \\ &\leq \exp \left( -\frac{\frac{1}{2}N^2p^2}{p(1-p) \sum_{i=1}^n c_i^2 + \frac{1}{3}MNp} \right) \\ &\leq \exp \left( -\frac{\frac{1}{2}N^2p^2}{pMN + \frac{1}{3}MNp} \right) = \exp \left( -\frac{3}{8} \frac{N}{M} p \right), \end{aligned}$$

where we used the fact that  $\sum_{i=1}^n c_i^2 \leq \sum_{i=1}^n c_i M = NM$ . ■

### 3 Two-Dimensional Balanced Allocations

In this section we analyze the maximal load of two schemes for the two-dimensional balanced allocation problem, as discussed in Section 1.3. Our first scheme is a one-choice process that is described and analyzed in Section 3.1, and our second scheme is a two-choice process that is described and analyzed in Section 3.2.

The input to the two-dimensional balanced allocation problem is a collection of lists,  $L_1, \dots, L_k$ , where the  $i$ th list  $L_i$  is of length  $n_i$ . The goal in this problem is to allocate an array of bins and to place each element from these lists in a bin (where a bin may contain several elements), such that given an index  $i$  we can efficiently recover the entire list  $L_i$ . Motivated by our goal of searchable symmetric encryption, we would like to design allocation scheme with good space overhead, locality, and read efficiency (while ignoring any security properties for now).

#### 3.1 A One-Choice Allocation Scheme

Our one-choice scheme is described as Algorithm 3.1. As discussed in Section 1.3, for each list  $L_i$  we choose a uniform bin, and then place the elements of  $L_i$  in consecutive bins starting from that bin (one element in each bin).

We show that for an appropriate choice of parameters, the maximal load is very close to its expectation. This is similar to the one-dimensional problem that considers balls instead of lists, and the locations of all balls are independent. In our case, the locations of the elements from the same list are clearly related, and the loads of any two consecutive bins are strongly correlated. This requires a more subtle analysis, and we prove the following theorem:

**Claim 3.2.** *Fix any  $m, k$ , and  $n_1, \dots, n_k$ . Let  $n = \sum_{i=1}^k n_i$  and assume that  $m \leq n$ . Then, with probability at least  $1 - m \cdot 2^{-n/m}$ , at the end of Algorithm 3.1 the maximal load of any bin is at most  $3n/m$ .*

**ALGORITHM 3.1** (Algorithm OneChoiceAllocation( $m, (n_1, \dots, n_k)$ )).

- **Input:** Number of bins  $m$ , and a vector of integers  $(n_1, \dots, n_k)$  representing the length of the lists  $L_1, \dots, L_k$ .
- **The algorithm:**
  1. Initialize  $m$  empty bins  $B_0, \dots, B_{m-1}$ .
  2. For each list  $L_i$  where  $i = 1, \dots, k$ :
    - (a) Sample  $\alpha \leftarrow \{0, 1, \dots, m-1\}$ .
    - (b) For  $j = 0, \dots, n_i - 1$ :
      - i. Place the  $j$ th element of the list  $L_i$  in bin  $B_{\alpha+j \bmod m}$ .

**Proof.** For simplicity, we first assume that there is no list with more than  $m$  elements. We later remove this assumption.

For  $0 \leq j \leq m-1$ , let  $X_j$  be a random variable denotes the load of bin  $B_j$ , and for every  $i = 1, \dots, k$ , let  $Y_j[i]$  be an indicator that gets 1 if and only if some element of the  $i$ th list falls into bin  $B_j$ . Note that  $X_j = \sum_{i=0}^{k-1} Y_j[i]$ . Moreover, for a fixed  $j \in \{0, \dots, m-1\}$ ,  $i \in \{1, \dots, k\}$ , we have:

$$E(Y_j[i]) = \frac{n_i}{m} .$$

This holds since there is no list with size greater than  $m$ , and therefore there is an element of some list  $L$  in bin  $B_j$  if and only if, the chosen placement for the head of the list is one of the values  $\{j, [(j-1) \bmod m], \dots, [(j-n_i+1) \bmod m]\}$ .

This implies that:

$$E[X_j] = E \left[ \sum_{i=0}^{k-1} Y_j[i] \right] = \sum_{i=0}^{k-1} E[Y_j[i]] = \frac{\sum_{i=0}^{k-1} n_i}{m} = \frac{n}{m} .$$

Although the random variables  $X_0, \dots, X_{m-1}$  are dependent (i.e., the loads of the bins), and even for every list the random variables  $Y_0[i], \dots, Y_{m-1}[i]$  are dependent (i.e., the bins where the elements of some single list are placed), for every fixed  $j \in \{0, \dots, m-1\}$  the random variables  $Y_j[1], \dots, Y_j[k]$  are *independent* (i.e., the choices of elements inside a given bin) . Also, they all taking values in  $\{0, 1\}$ . Therefore, we can apply Chernoff's bound (Lemma 2.10), and get:

$$\Pr [X_j \geq (1 + \delta) E[X_j]] \leq \left( \frac{e^\delta}{(1 + \delta)^{(1+\delta)}} \right)^{E[X_j]}$$

For  $\delta \geq e - 1$ , we have that:

$$\frac{e^\delta}{(1 + \delta)^{(1+\delta)}} = \frac{1}{1 + \delta} \cdot \left( \frac{e}{1 + \delta} \right)^\delta \leq \frac{1}{1 + \delta} \leq \frac{1}{2} ,$$

and thus

$$\Pr \left[ X_i \geq 3 \cdot \frac{n}{m} \right] \leq 2^{-n/m} .$$

Using union-bound, we conclude that the probability that one of the bins has load greater than  $3 \cdot n/m$  is at most  $m \cdot 2^{-n/m}$ .

For case where there exist lists with size greater than  $m$ , observe that for each such a list  $i$ , each bin will have at least  $\lfloor n_i/m \rfloor$  element of that list, and exactly  $\lfloor n_i \bmod m \rfloor$  bins will have an additional one element from that list. For each such a list and for each bin  $B_i$ , instead of considering the single



random variable  $Y_j[i]$ , we define  $\lfloor n_i/m \rfloor + 1$  random variables  $Y_j^{(0)}[i], \dots, Y_j^{(\lfloor n_i/m \rfloor)}[i]$ , where  $Y_j^{(\gamma)}[i]$  indicates whether at least  $\gamma$  elements of the  $i$ th list fell into bin  $B_j$ . The first  $\lfloor n_i/m \rfloor$  variables are constant, and always get 1. The random variable  $Y_j^{(\lfloor n_i/m \rfloor)}[i]$  is 1 with probability  $\lfloor n_i \bmod m \rfloor / m$ . The expected sum of these random variables is  $n_i/m$ , and therefore  $E[X_j] = n/m$ , as in the simplified case. Again, these variables are independent and get values in  $\{0, 1\}$ , so we can apply Chernoff's bound and get exactly the same result. ■

When constructing our symmetric searchable encryption schemes, we are interested in having a negligible failure probability, and we derive the following corollary from Claim 3.2:

**Corollary 3.3.** *Fix any  $k$ , and  $n_1, \dots, n_k$ . Let  $n = \sum_{i=1}^k n_i$  and  $m = \frac{n}{\log n \cdot \log \log n}$ . Then, with probability at least  $1 - n^{-\omega(1)}$ , at the end of Algorithm 3.1 the maximal load of any bin is at most  $3 \log n \cdot \log \log n$ .*

### 3.2 A Two-Choice Allocation Scheme

Our next step is generalizing our one-choice allocation scheme into a two-choice one. Before describing our scheme, we first discuss several possible generalizations that seem very reasonable and natural, but turn out somewhat insufficient and do not provide the desired properties.

Consider an algorithm that for every list  $L_i$  chooses two possible bins,  $B_{\alpha_1}$  and  $B_{\alpha_2}$ , and places the elements of the list in a consecutive manner starting from the least loaded bin  $B_{\alpha_{i^*}}$  among these two choices. However, this approach seems insufficient, since the fact that  $B_{\alpha_{i^*}}$  is the least loaded bin among the two possible bins  $B_{\alpha_1}$  and  $B_{\alpha_2}$  does not imply that  $B_{\alpha_{i^*+j}}$  is the least loaded bin among  $B_{\alpha_1+j}$  and  $B_{\alpha_2+j}$  for any  $j > 1$ . That is, the decision we make in order to place the head of the list does not necessarily apply for the placement of the following element of the list.

Another natural approach is the following. Similarly to the previous proposal, for each list we choose two possible locations,  $B_{\alpha_1}$  and  $B_{\alpha_2}$ , and place its first element in the least loaded bin among these two possibilities. However, unlike the previous proposal, we make a similar decision for each element of the list. Specifically, the  $j$ th element of the list is placed in the least loaded bin among  $B_{\alpha_1+j-1}$  and  $B_{\alpha_2+j-1}$ . Unfortunately, this algorithm seems somewhat difficult to analyze as it introduces many (redundant) dependencies between elements, and there are too many decisions that are made in the algorithm. A more elegant approach would be to have a single comparison for the placement of the entire list.

We now turn to describe our two-choice approach. Assume without loss of generality that the lengths of all lists are powers of two (otherwise we can pad each with dummy elements and thus increase the database size at most by a factor of 2), and that the number of bins is also a power of two. In our algorithm, we process the lists according to their length. We first sort the lists in a descending order, and then, for a given list of size  $n_i$ , we view each block of  $n_i$  consecutive bins as a single possible “super bin”. As a result, the possible locations for the head of the list is not  $m$  (the overall number of bins), but rather only  $m/n_i$ , the number of “super bins”. Then, we choose two possible super bins, and place the entire list, in the least loaded super bin among these two possible choices. The actual placement is done element by element, in the respective bins that constitute this super bin (i.e., one element per “standard” bin).

This approach has the following advantages: First, the entire list is placed according to a *single* comparison, which significantly simplifies the analysis. Second, all the elements of each list are placed *at the exact same level*. That is, our algorithm has the following invariant: Whenever we place a list, all the bins that constitute a possible super bin have the exact same number of elements. This is mainly due to the fact that we place the lists in descending order, and in each iteration the

super bins are split into smaller super bins that each matches the size of the current list. A formal description is presented as Algorithm 3.4.

**ALGORITHM 3.4** (Algorithm TwoChoiceAllocation( $m, (n_1, \dots, n_k)$ )).

- **Input:** Number of bins  $m$  and a vector of integers  $(n_1, \dots, n_k)$  representing the length of the lists  $L_1, \dots, L_k$ . We let  $n = \sum_{i=1}^k n_i$ , and assume for concreteness that  $m$  and the  $n_i$ 's are powers of two, and that  $m \geq n_1 \geq n_1 \geq \dots \geq n_k$ .
- **The algorithm:**
  1. Initialize  $m$  empty bins  $B_0, \dots, B_{m-1}$ .
  2. For each list  $L_i$  where  $i = 1, \dots, k$ :
    - (a) Choose  $\alpha_1, \alpha_2 \leftarrow \{0, \dots, \frac{m}{n_i} - 1\}$  independently and uniformly at random. Consider the two super bins  $\tilde{B}_{\alpha_1} = (B_{n_i \cdot \alpha_1 + j})_{j=0}^{n_i-1}$  and  $\tilde{B}_{\alpha_2} = (B_{n_i \cdot \alpha_2 + j})_{j=0}^{n_i-1}$ .
    - (b) Let  $\tilde{B}_\alpha$  be the least loaded<sup>7</sup> among  $\tilde{B}_{\alpha_1}$  and  $\tilde{B}_{\alpha_2}$ . We place the list  $L_i$  in the super bin  $\tilde{B}_\alpha$ . That is, for every  $j = 0, \dots, n_i - 1$ :
      - i. Place the  $j$ th element of the list  $L_i$  in the bin  $B_{n_i \cdot \alpha + j}$ .

### 3.2.1 Bounding the Maximal Load

We prove the following theorem:

**Theorem 3.5.** *Assume that each  $n_i$  is a power of two and that  $n_1 \geq \dots \geq n_k$ . Let  $S \geq n_1$  be some bound on the lengths, and let  $m$  be the number of bins. Then*

1. *With probability  $1 - n^{-\Omega(\log n)} - e^{-\Omega(m^3/(n^2 \cdot S))}$ , there are at most  $S \log^2 n$  elements at level greater than  $\frac{8n}{m} + \log \log \frac{n}{S} + 2$ .*
2. *Let  $\epsilon > 0$ ,  $S = n^{1-\epsilon}$ , and assume that  $m \geq \frac{n}{\log n}$ . Then, for any non-decreasing function  $f(n)$  satisfying  $f(n) = \Omega(\log \log n)$ ,  $f(n) = O(\sqrt{\log n})$  and  $f(2n) = O(f(n))$ , the maximal load is at most  $\frac{8n}{m} + O(\log \epsilon^{-1} \cdot f(n))$  with probability  $1 - O(\log \epsilon^{-1}) \cdot n^{-\Omega(\epsilon \cdot f(n^\epsilon))} - e^{-\Omega(m^3/(n^2 \cdot S))}$ .*

For our searchable symmetric encryption schemes, we will rely on the following corollary (say, with  $A = 1$ ) that follows from Theorem 3.5:

**Corollary 3.6.** *For any constant  $A \geq 1$ , let  $\epsilon = (\log \log n)^{-A}$  (and so the size of the maximal list is bounded by  $n^{1-1/(\log \log n)^A}$ ) and  $m = \frac{n}{(\log \log n)^A \cdot (\log \log \log n)^2}$ . Then, with overwhelming probability the maximal load is  $O((\log \log n)^A \cdot (\log \log \log n)^2)$ .*

**Proof overview.** The proof is rather subtle, and is inspired by the “layered induction” technique of the one-dimensional problem [ABK<sup>+</sup>99] (see also [MRS00, MU05]). The following overview is rather simplified, and as a result some of its notation is not identical to that of the formal proof below.

In order to bound the maximal load, we need to approximately bound the number of bins with  $j$  elements for every value of  $j$ . The proof defines a sequence of values  $\beta_1, \beta_2, \dots$  and shows that the number of bins with load at least  $j$  is bounded by  $\beta_j$  with high probability. The proof follows by induction, where in each step, we show that if the number of bins with load at least  $j$  is at most

<sup>7</sup>The load of a super bin is defined as the sum of loads of the “standard” bins that constitutes that super bin. According to the invariant of this algorithm, the load of these standard bins is equal, and therefore the sum is simply the number of standard bins that constitute the super bin times the load of each one of them.

$\beta_j$ , then with high probability, the number of bins with load at least  $j + 1$  is at most  $\beta_{j+1}$ . The values  $\beta_1, \beta_2, \dots$  are defined according to a recursive relation. Intuitively, the sequence  $\beta_1, \beta_2, \dots$  is monotonic decreasing, and we look for the first  $j$  for which  $\beta_j < 1$  (i.e., there are no bins with load at least  $j$ ). As we will see, this description is a bit simplified and we cannot actually use the recursive relation all the way to  $\beta_j < 1$ , since the induction does not hold unconditionally. Thus, we will have to stop at a larger  $\beta_j$  and conclude the proof using some additional arguments. We will elaborate on this point later on.

**The recursive relation.** According to our algorithm, when we place the  $i$ th list, all of its elements share the exact same level. Therefore, a list will have height (or, level) at least  $j + 1$  only if the two choices for its super bins are of level at least  $j$ . Since there are at most  $\beta_j$  bins with load at least  $j$  (or,  $\beta_j/n_i$  super bins of size  $n_i$ ) and  $m$  bins (or,  $m/n_i$  super bins of size  $n_i$ ), the probability that the two possible choices have level at least  $j$  is  $(\beta_j/m)^2$ . Therefore, the expected number of list elements of height at least  $j + 1$  is at most  $n \cdot (\beta_j/m)^2$ , and so the number of list elements of height at least  $j + 1$  is no more than twice this amount with high probability (using some tail bound). This implies that there are no more than  $2n \cdot (\beta_j/m)^2$  bins with height at least  $j + 1$  and this leads to the following recursive relation:

$$\beta_{j+1} \stackrel{\text{def}}{=} 2n \cdot \left(\frac{\beta_j}{m}\right)^2.$$

**The induction.** We confirm the recursive relation by induction. We prove that with high probability, the number of bins with list elements of height at least  $j + 1$  is bounded by  $\beta_{j+1}$ , conditioned that the number of bins with load at least  $j$  is bounded by  $\beta_j$ . This step is somewhat subtle in the original proof, primarily because one must handle the conditioning appropriately, and the formal argument requires some care. In our case, the argument is even more subtle. In the original proof, this step is shown using Chernoff's bound, however, we cannot use this bound here due to the fact that the random variables do not take values in  $\{0, 1\}$ , and depend on the lengths of the lists. We use Bernstein's bound instead, and this step already introduces our limitation on the length of the longest list.

**Concluding the first part of Theorem 3.5.** The proof then proceeds as follows, where for this high-level overview, we assume for simplicity that  $m = n$ . We may set  $\beta_4 = n/4$ , since there cannot be more than  $n/4$  bins with at least 4 elements. As in the original proof, the induction step does not hold unconditionally, and therefore we cannot use the recursion all the way to the first index  $k$  for which  $\beta_k < 1$ . Instead, we have to stop earlier, when  $\beta_j = S \log^2 n$  where  $S$  is some upper bound on the length of the longest list. It turns out that the original proof is a special case of ours, where this layered induction holds as long as  $\beta_j \geq \log^2 n$  (i.e.,  $S = 1$ ). The first index  $j^*$  for which  $\beta_{j^*} = S \log^2 n$  is when  $j^* = O(\log \log(n/S))$ . This means that there are no more than  $S \log^2 n$  list elements at level greater than  $O(\log \log(n/S))$ , proving the first part of Theorem 3.5.

**The second part of Theorem 3.5.** At this stage, we cannot use the tail bound any more and we need a different (more coarse) argument. Clearly, if there is a bin with load at least  $j^* + a$ , then there are at least  $a$  lists with height at least  $j^* + 1$ . Assuming that there are no more than  $\beta_{j^*} = S \log^2 n$  bins with load at least  $j^*$ , using a simple union bound argument, we can bound the probability that there are at least  $a$  lists with height at least  $j^* + 1$  by roughly  $(S^2/n)^a$ . This is negligible if  $a = O(\log \log n)$  and  $S = n^\delta$  for some constant  $\delta < 1/2$ . Combining this with the above,

we conclude that the maximal load of the process is  $O(\log \log n)$  assuming that there are no lists with more than  $n^\delta$  elements.

**A scaling argument.** The above shows that the algorithm places lists of length up to  $n^\delta$  for some constant  $\delta < 1/2$  (say, for concreteness,  $\delta = 1/3$ ). We claim that the algorithm actually succeeds in placing much longer lists, and we show this using a more subtle argument.

Our first observation is that the algorithm can be “scaled”. Namely, consider the following two executions of the algorithm. In the first execution, we have list of lengths  $(n_1, \dots, n_k)$  and number of bins  $m$  where all are multiples of some  $s$ . In the second execution, instead of running the algorithm on these inputs, we divide the elements to blocks of size  $s$ , run the algorithm on the input  $(n_1/s, \dots, n_k/s)$  and  $m/s$ , and then extend each bin to  $s$  bins (by placing the blocks in the respective bins). This, in fact, is what implicitly happens in the first execution, and therefore these two executions are equivalent.

As a matter of fact, using the behavior of the algorithm on the latter case we can conclude its properties on the former. Specifically, let  $n' = n/s$  (i.e.,  $n'$  is the sum of lengths in the scaled case). The failure probability of this execution is about  $(S^2/n')^a$  and it can handle lists of length at most  $n'^{1/3}$ , which implies that the former case in fact fails with probability about  $(S^2/(n/s))^a$  and handles lists of lengths at most  $(n/s)^{1/3}$ . By taking  $s = n^{1/3}$ , we obtain that the algorithm can handle lists of lengths between  $s = n^{1-(2/3)^1}$  and  $n^{1-(2/3)^2}$  with all but negligible probability. We can repeat this argument carefully, and show that the algorithm can handle lists of length between  $n^{1-(2/3)^\alpha}$  and  $n^{1-(2/3)^{\alpha+1}}$  with all but negligible probability, for any constant  $\alpha > 0$ .

Our second observation is that the algorithm can in fact handle all of these possible inputs *simultaneously*. This is because the algorithm behaves similarly to an algorithm that divides the lists into sets of inputs of lengths between  $n^{1-(2/3)^\alpha}$  and  $n^{1-(2/3)^{\alpha+1}}$ , invokes the algorithm on each set of inputs independently (i.e., on a different set of bins), and then combines the bins of all these executions. By carefully defining the number of sets, we conclude that the maximal load of the latter algorithm is at most  $\tilde{O}(\log \log n)$ , implying that our algorithm behaves similarly to that as well.

**Proof of Theorem 3.5.** We will do our calculations with respect to some bound  $N \geq n$ . We point out an important loop invariant: Before the loop  $i$ , for each  $0 \leq \beta \leq \frac{m}{n_i} - 1$ , the current load of all the bins in the super bin  $\tilde{B}_\beta = (B_{n_i, \beta+j})_{j=0}^{n_i-1}$  is the same. So all the elements of the list  $L_i$  are placed at the same height. We refer to that common height as the height of the list  $L_i$ , and denote it by  $h(i)$ .

We make further notations: Let  $\ell = \lceil \frac{8n}{m} \rceil$ .  $X_i$  is where we placed the head of the list  $L_i$ . We refer to the stage right after placing the list  $L_i$  as “time  $i$ ”.  $\nu_j(i)$  is the number of bins with load  $\geq \ell + j$  at time  $i$ .  $\mu_j(i)$  is the number of list *elements* with height  $\geq \ell + j$  at time  $i$  (mathematically,  $\sum_{s \leq i \wedge h(s) \geq \ell+j} n_s$ ). We also define  $\nu_j(0) = \mu_j(0) = 0$ . Trivially we have  $\nu_j(i) \leq \mu_j(i)$ . We define  $\beta_0 = \frac{m^2}{4n}$ ,  $\beta_{j+1} = 2n \frac{\beta_j^2}{m^2}$  and  $\mathcal{E}_j = \{\nu_j(k) \leq \beta_j\}$  (i.e.,  $\mathcal{E}_j$  is the event that in the end of the algorithm’s execution there are at most  $\beta_j$  bins with load  $\geq \ell + j$ ). We claim the following:

**Claim 3.7.** *There exists a constant  $c > 0$  such that*

$$\Pr[-\mathcal{E}_0] \leq \exp\left(-\frac{c \cdot m^3}{n^2 \cdot S}\right)$$

We assume the correctness of Claim 3.7 for now and prove it later. Fix a level  $\ell + j$ . For  $i \in [k]$

define the random variable  $Y_i$  as follows,

$$Y_i = \begin{cases} 1 & h(i) \geq \ell + j + 1 \wedge \nu_j(i-1) \leq \beta_j \\ 0 & \text{else} \end{cases}$$

Conditioned on  $\mathcal{E}_j$  we have  $\mu_{j+1}(k) = \sum_{i=1}^k n_i Y_i$ . Letting  $p_j = \frac{\beta_j^2}{m^2}$ , we have  $Y_i = Y_i(X_1, \dots, X_i)$  and  $\Pr[Y_i = 1 | X_1, \dots, X_{i-1}] \leq p_j$ , so by applying Corollary 2.14 we get

$$\Pr \left[ \sum_{i=1}^k n_i Y_i > \beta_{j+1} \right] = \Pr \left[ \sum_{i=1}^k n_i Y_i > 2np_j \right] \leq \exp \left( -\frac{3n}{8S} p_j \right) \leq N^{-c \log N}$$

where the last step holds as long as  $\frac{2n}{S} p_j \geq \log^2 N$ , therefore

$$\begin{aligned} \Pr[\neg \mathcal{E}_{j+1}] &\leq \Pr[\neg \mathcal{E}_{j+1} \wedge \mathcal{E}_j] + \Pr[\neg \mathcal{E}_j] \\ &\leq \Pr[\mu_{j+1}(k) > \beta_{j+1} \wedge \mathcal{E}_j] + \Pr[\neg \mathcal{E}_j] \\ &\leq \Pr \left[ \sum_{i=1}^k n_i Y_i > \beta_{j+1} \right] + \Pr[\neg \mathcal{E}_j] \\ &\leq N^{-c \log N} + \Pr[\neg \mathcal{E}_j] \end{aligned}$$

Let  $j^*$  be the minimal  $j$  such that  $p_j < \frac{S}{2n} \log^2 N$ . We shall estimate  $j^*$ . By induction one can see that

$$\beta_j \leq \frac{m^2}{n \cdot 2^{2j}}, \quad p_j = \frac{\beta_j^2}{m^2} \leq \left( \frac{m}{n \cdot 2^{2j}} \right)^2.$$

We have

$$\left( \frac{m}{n \cdot 2^{2j^*-1}} \right)^2 \geq p_{j^*-1} \geq \frac{S}{2n} \log^2 N,$$

and by applying log we get

$$2j^* \leq \log \frac{n}{S} - 2 \log \frac{n}{m} - 2 \log \log N + 1 \leq \log \frac{n}{S}$$

so  $j^* \leq \log \log \frac{n}{S}$ . Now we apply the bound one more time with  $p'_{j^*} = \frac{S}{2n} \log^2 N > p_{j^*}$ . We set  $\beta'_{j^*+1} = 2np'_{j^*} = S \log^2 N$ ,  $\mathcal{F} = \{\mu_{j^*+1}(k) \leq \beta'_{j^*+1}\}$  and  $\mathcal{E}'_{j^*+1} = \{\nu_{j^*+1}(k) \leq \beta'_{j^*+1}\}$ . As before we get  $\Pr[\neg \mathcal{F}] \leq N^{-c \log N} + \Pr(\neg \mathcal{E}_{j^*})$ , so we conclude by induction that  $\Pr(\neg \mathcal{E}'_{j^*+1}) \leq \Pr[\neg \mathcal{F}] \leq (j^* + 1) \cdot N^{-c \log N} + \Pr[\neg \mathcal{E}_0] \leq N^{-\Omega(\log N)} + \Pr[\neg \mathcal{E}_0]$ , so

$$\Pr[\mathcal{F}] = \Pr[\mu_{j^*+1}(k) \leq S \log^2 N] = 1 - N^{-\Omega(\log N)} - \Pr[\neg \mathcal{E}_0],$$

and by setting  $N = n$  and using Claim 3.7 we deduce the first item of Theorem 3.5.

Now, as a step towards proving the second item, we temporarily assume that  $m \geq \frac{N}{\log^2 N}$  and  $S = N^\delta$  for some  $0 \leq \delta < \frac{1}{2}$ . For the level  $\ell + j^* + 1$  we simply define

$$Y'_i = \begin{cases} 1 & h(i) \geq \ell + j^* + 2 \\ 0 & \text{else} \end{cases}$$

then we have

$$\Pr[Y'_i = 1 \wedge \mathcal{E}'_{j^*+1}] \leq \frac{\beta'^2_{j^*+1}}{m^2} = \frac{S^2}{m^2} \log^4 N \leq \frac{\log^8 N}{N^{2-2\delta}}.$$

If  $\nu_{j^*+a+1}(k) \geq 1$  then there at least  $a$  lists with height  $\geq \ell + j^* + 2$ , therefore by a union bound we get

$$\begin{aligned} \Pr[\nu_{j^*+a+1}(k) \geq 1 \wedge \mathcal{E}'_{j^*+1}] &\leq \Pr\left[\sum_{i=1}^k Y'_i \geq a \wedge \mathcal{E}'_{j^*+1}\right] \\ &\leq \binom{k}{a} \left(\frac{\log^8 N}{N^{2-2\delta}}\right)^a \\ &\leq N^a \left(\frac{\log^8 N}{N^{2-2\delta}}\right)^a = \left(\frac{\log^8 N}{N^{1-2\delta}}\right)^a \end{aligned}$$

so by taking  $a = f(N)$  we get

$$\begin{aligned} \Pr[\nu_{j^*+a+1}(k) \geq 1] &\leq \Pr[\nu_{j^*+a+1}(k) \geq 1 \wedge \mathcal{E}'_{j^*+1}] + \Pr[\neg \mathcal{E}'_{j^*+1}] \\ &\leq \left(\frac{\log^8 N}{N^{1-2\delta}}\right)^{f(N)} + N^{-\Omega(\log N)} + \Pr[\neg \mathcal{E}_0] \\ &= N^{-\Omega(f(N))} + \Pr[\neg \mathcal{E}_0] \end{aligned}$$

for the level  $\ell + j^* + a + 1 = \frac{8n}{m} + O(f(N))$ . Note that the constants here depend on  $\delta$ , but we will actually consider a specific case, namely  $\delta = \frac{1}{3}$ . Now, we prove the second item of Theorem 3.5 by the following scaling argument. First we point out some observations:

- Our algorithm is scalable, i.e., if every list is larger than some power of two  $s$ , and we define a new input by dividing the length of every list,  $n'_i = \frac{n_i}{s}$ , and also the number of bins is  $m' = \frac{m}{s}$ , then the algorithm works on the original input as if it worked on the new input and scaled the result by duplicating each bin  $s$  times.
- Suppose  $S = N^{\frac{1}{3}}$ ,  $m \geq \frac{N}{\log^2 N}$ . Then, we saw that the maximal load exceeds  $\frac{8n}{m} + O(f(N))$  with probability at most  $N^{-\Omega(f(N))} + \Pr[\neg \mathcal{E}_0]$ . Moreover, the same argument shows that if we already placed some other larger lists such that the maximal load is at most  $\ell + t$  for an integer  $t \geq 0$ , then as long that the event  $\mathcal{E}_0$  occurs, we will increase the current maximal load to at least  $\ell + t + O(f(N))$  with probability at most  $N^{-\Omega(f(N))}$ .

Now, we assume that  $S = n^{1-\epsilon}$  and  $m \geq \frac{n}{\log n}$  as in the second item of Theorem 1.1. We may assume that  $\epsilon \geq \frac{1}{\sqrt{\log n}}$  since the claim is vacuously true for smaller  $\epsilon$ . We set  $\gamma = \frac{2}{3}$  and let  $z = \lceil \log_\gamma \epsilon \rceil$ . Then  $\gamma^z \leq \epsilon < \gamma^{z-1}$ . We divide the lists into  $z$  sets by  $V_j = \{i \mid n^{1-\gamma^j} \leq n_i \leq n^{1-\gamma^{j+1}}\}$  for  $j = 0, \dots, z-1$ . We scale the lists in  $V_j$  by  $s = n^{1-\gamma^j}$ ,  $n'_i = \frac{n_i}{s}$ . Also, we scale then number of bins  $m' = \frac{m}{s}$ . Now the input size is  $n' = \sum_{i \in V_j} n'_i$ . Let  $N' = \frac{n}{s} \geq n'$ . We have  $N' = n^{\gamma^j} \geq n^\epsilon$ , so

$$m' = \frac{N'}{\log n} = \frac{N'}{\log^2 N'} \cdot \frac{\log^2 N'}{\log n} \geq \frac{N'}{\log^2 N'} \cdot \epsilon^2 \log n \geq \frac{N'}{\log^2 N'}$$

Each list is of length at least  $\frac{n^{1-\gamma^{j+1}}}{n^{1-\gamma^j}} = n^{\gamma^j(1-\gamma)} = (N')^{\frac{1}{3}}$ , So by the scaling observation, assuming  $\mathcal{E}_0$  occurs, and that the maximal load before inserting the lists in  $V_j$  is at most  $\ell + t$  for an integer  $t \geq 0$ , the maximal load after inserting the lists in  $V_j$  will be at least  $\ell + t + O(f(N'))$  with probability at most  $(N')^{-\Omega(f(N'))} \leq n^{-\Omega(\epsilon \cdot f(n^\epsilon))}$ . So by a union bound, we get that the maximal load is at most

$$\ell + \sum_{i=0}^{z-1} O(f(n^{\gamma^i})) \leq \frac{8n}{m} + O(\log \epsilon^{-1} \cdot f(n))$$

with probability  $1 - O(\log \epsilon^{-1}) \cdot n^{-\Omega(\epsilon \cdot f(n^\epsilon))} - \Pr[-\mathcal{E}_0]$ , and using Claim 3.7 we deduce the second item of Theorem 3.5.  $\blacksquare$

It remains to prove Claim 3.7.

**Proof of Claim 3.7.** We remind that  $\mathcal{E}_0 = \{\nu_0(k) \leq \beta_0\}$ , where  $\nu_0(k)$  is the number of bins with load  $\geq \ell = \lceil \frac{8n}{m} \rceil$  at the end of the algorithm's execution,  $\beta_0 = \frac{m^2}{4n}$ , and we are interested in bounding  $\Pr[\mathcal{E}_0]$  for above. If  $n < 2m$ , then there can be at most  $n/\ell$  bins with load  $\geq \ell$ , but  $n/\ell < 2m/(8n/m) = \beta_0$ , so  $\Pr[-\mathcal{E}_0] = 0$  and the claim follows.

From now on we may assume that  $n \geq 2m$ . Consider the allocation algorithm in which we place each list  $L_i$  in both sampled super bins. Clearly, we can bound  $\Pr[-\mathcal{E}_0] \leq \Pr[\nu'_0 > \beta_0]$ , where  $\nu'_0$  is the number of bins with load  $\geq \ell$  at the end of this new algorithm's execution. To the end of this section, all of our discussion is with respect to this new algorithm.

For each  $0 \leq i \leq m-1$ , we define the random variable

$$X_i = \begin{cases} 1 & \text{the load of } B_i \text{ is larger than } \beta_0 \\ 0 & \text{else} \end{cases}.$$

We remind that  $S$  upper bounds the length of the maximal list, and we may assume w.l.o.g. that  $S$  is a power of 2. For each  $0 \leq j \leq m/S - 1$ , we define  $Y_j = \sum_{i=(j-1) \cdot S}^{j \cdot S - 1} X_i$ . Clearly  $\sum_i X_i = \sum_j Y_j$ , where the sums are with respect to the indices for which we defined the variables. A proof similar to the one-choice case shows that  $\Pr[X_i = 1] \leq 2^{-2n/m}$  (i.e., representing the load of the bin  $B_i$  as a sum of indicators and applying Chernoff's bound). Therefore,  $\mathbb{E} \sum_j Y_j = \mathbb{E} \sum_i X_i \leq m \cdot 2^{-2n/m}$ . For every  $x \geq 2$  it holds that  $2^{2x} \geq 8x$ , so since  $n/m \geq 2$ , it holds that  $2^{-2n/m} \leq m/8n$ , and therefore  $\mathbb{E} \sum_j Y_j \leq \beta_0/2$ . We want to use Hoeffding's bound on the variables  $Y_0, \dots, Y_{m/S-1}$ , so we prove the following claim.

**Claim 3.8.** *The variables  $Y_0, \dots, Y_{m/S-1}$  are negatively associated.*

**Proof.** For integers  $1 \leq i \leq k$ ,  $b \in \{1, 2\}$  and  $0 \leq \alpha \leq m/n_i - 1$ , we define the random variable

$$Z_{i,b,\alpha} = \begin{cases} 1 & \text{the } b\text{th super bin sampled for the list } L_i \text{ is } \tilde{B}_\alpha = (B_{n_i \cdot \alpha + j})_{j=0}^{n_i-1} \\ 0 & \text{else} \end{cases}.$$

Fixing  $i$  and  $b$ , it holds that  $\sum_{\alpha=0}^{m/n_i-1} Z_{i,b,\alpha} = 1$ , and therefore by the zero-one principle (Lemma 2.7), the variables  $Z_{i,b,0}, \dots, Z_{i,b,m/n_i-1}$  are negatively associated. For all  $i$  and  $b$ , let  $W_{i,b} = (Z_{i,b,0}, \dots, Z_{i,b,m/n_i-1})$ . Then, the vectors  $W_{1,0}, W_{1,1}, \dots, W_{k,0}, W_{k,1}$  are independent and each of them is a vector of negatively-associated variables, therefore, by Lemma 2.8 their concatenation  $(Z_{i,b,\alpha})_{i,b,\alpha}$  is a vector of negatively-associated variables. For each  $0 \leq j \leq m/S - 1$ , the variable  $Y_j$  is a monotone function of the variables  $(Z_{i,b,\alpha})_{(i,b,\alpha) \in V_j}$ , where

$$V_j = \{(i, b, \alpha) \mid (j-1) \cdot S \leq \alpha \cdot n_i \leq j \cdot S - 1\}.$$

The sets  $V_0, \dots, V_{m/S-1}$  are pairwise disjoint, and therefore by Lemma 2.9, the variables  $Y_0, \dots, Y_{m/S-1}$  are negatively associated. This concludes the proof of Claim 3.8.  $\blacksquare$

Now, we can indeed apply Hoeffding's bound (Lemma 2.11). For each  $0 \leq j \leq m/S - 1$  it holds that  $Y_j \in [0, S]$ , so the sum term in Hoeffding's bound calculates to  $\sum_{j=0}^{m/S-1} (b_j - a_j)^2 =$

$\sum_{j=0}^{m/S-1} S^2 = m \cdot S$ , and we obtain that

$$\begin{aligned} \Pr \left[ \sum_{i=0}^m X_i > \beta_0 \right] &= \Pr \left[ \sum_{j=0}^{m/S-1} Y_j > \beta_0 \right] \\ &\leq \Pr \left[ \sum_{j=0}^{m/S-1} Y_j - \mathbb{E} \sum_{j=0}^{m/S-1} Y_j > \frac{\beta_0}{2} \right] \\ &\leq \exp \left( -\frac{2 \cdot (\beta_0/2)^2}{m \cdot S} \right) \\ &= \exp \left( -\frac{m^3}{32 \cdot n^2 \cdot S} \right), \end{aligned}$$

and this settles the proof of Claim 3.7.  $\blacksquare$

### 3.2.2 The Tightness of Our Analysis

The above analysis assumes a bound on the length of the longest list. Specifically, we assumed an upper bound  $n^{1-\epsilon(n)}$  on the length of the longest list, where  $\epsilon = (\log \log n)^{-A}$ , and proved an upper bound  $g(n)$  on the maximal load of Algorithm 3.4 with an overwhelming probability, where  $g(n) = \tilde{O}((\log \log n)^A)$ . The following claim shows that this trade-off is in fact tight, namely that  $g(n) = \omega((\log \log n)^A)$ :

**Claim 3.9.** *For any constants  $A \geq 0$  and  $c \geq 0$ , and for any sufficiently large enough  $n$ , there exists an input for Algorithm 3.4 such that (1) the number of bins is  $m = n$ , (2) the length of the longest list is  $n^{1-\epsilon}$  where  $\epsilon = (\log \log n)^{-A}$ , and (3) the maximal load is at least  $c \cdot (\log \log n)^A$  with a noticeable probability.*

**Proof.** Consider the following input to Algorithm 3.4. The number of bins is  $m = n$ , and the set of lists consist of  $L \leq n^\epsilon$  lists of the length  $n^{\stackrel{\text{def}}{=}n^{1-\epsilon}}$ , and an arbitrary number of lists (and structure) for the remaining  $n - L \cdot n^{1-\epsilon}$  elements (but still with no list longer than  $n'$ ).

We compute the probability that all two choices of all the  $L$  long lists share the exact same super bin (therefore, the load of each of the bins that constitute that super bin is at least  $L$ ). Since we have independence between different lists, this probability is

$$\left( \frac{n'}{m} \right)^{2 \cdot L} = \left( \frac{n^{1-\epsilon}}{n} \right)^{2 \cdot L} = n^{-2L\epsilon}.$$

Taking  $\epsilon = (\log \log n)^{-A}$  and  $L = c \cdot (\log \log n)^A$ , the above occurs with a noticeable probability and the maximal load is at least  $c \cdot (\log \log n)^A$ . Note that  $n$  should be large enough such that  $L \leq n^\epsilon$ .  $\blacksquare$

We note that the above argument can be generalized for a similar  $d$ -choice process for  $d > 2$ , and for a number of bins which is as large as  $n \log n$ . In addition, in Section 4.4 we generalize the above to a more general family of allocation algorithms.

## 4 SSE Schemes from Two-Dimensional Balanced Allocations

In this section we present a general framework for basing searchable symmetric encryption schemes on two-dimensional balanced allocations algorithms. First, in Section 4.1 we formally define the notion of an allocation algorithm, and derive concrete instantiations based on our two-dimensional allocation schemes from Section 3. Then, in Section 4.2 we show how to transform any allocation algorithm into a searchable symmetric encryption scheme using cryptographic techniques.



## 4.1 Allocation Algorithms

An allocation algorithm receives as input a database  $DB = \{DB(w_1), \dots, DB(w_{n_w})\}$  and places its elements in an array. For each list  $DB(w_i)$ , we distinguish between its *possible* locations in the array, and its *actual* locations in the array. We restrict our attention to allocation algorithms in which the *possible* locations of each list  $DB(w_i)$  depend only on its length  $n_i = |DB(w_i)|$  and on the size  $N = \sum_{i=1}^{n_w} |DB(w_i)|$  of the database. In particular, the possible locations of different lists are independent (in contrast, the actual locations of different lists are naturally allowed to be dependent). We model this property by considering allocation algorithms that follow a two-step structure. First, they separately generate the possible locations for each list using a procedure denoted `RangesGen`. Then, given the entire collection of possible locations of all list, they determine the actual locations of the lists using a procedure denote `Allocation`. The structure of the algorithm is described as Structure 4.1.

**STRUCTURE 4.1** (The structure of the allocation algorithm `Allocation-Alg`).

- **Input:**  $k$  integer values  $(n_1, \dots, n_k)$  representing the lengths of the lists. Let  $N = \sum_{i=1}^k n_i$ .
- **The algorithm:**
  1. For every  $1 \leq i \leq k$  let  $R_i \leftarrow \text{RangesGen}(N, n_i)$ .  
(The procedure `RangesGen`( $N, n_i$ ) outputs the possible ranges  $R_i = \{[a_1, b_1], \dots, [a_d, b_d]\}$  for the  $i$ th list.)
  2. Let  $\text{map} \leftarrow \text{Allocation}((n_1, R_1), \dots, (n_k, R_k))$ .  
(The array `map` holds the actual locations of the lists. Each entry in this array contains either a pair  $(i, j)$  representing that this entry is the actual location of the  $j$ th element from the  $i$ th list, or `NULL` representing an empty entry.)
  3. Output `map`.

We sometimes find it convenient to denote by `RangesGen`( $N, n_i; r$ ) an invocation of the procedure `RangesGen` on input  $(N, n_i)$  using specific randomness  $r$ .

**Efficiency measures.** We measure the efficiency of an allocation algorithm of the above structure with respect to its space overhead, its locality, and its read efficiency.

**Definition 4.2.** We say that `Allocation-Alg` = (`RangesGen`, `Allocation`) is an  $(s, d, r)$ -allocation algorithm, for some functions  $s(\cdot)$ ,  $d(\cdot)$  and  $r(\cdot)$ , if the following properties hold:

- **Correctness:** There exists a negligible function  $\text{negl}(\cdot)$  such that for every input  $(n_1, \dots, n_k)$

$$\Pr[(\text{Allocation-Alg}(n_1, \dots, n_k) = \perp) \vee \text{invalid-allocation}] \leq \text{negl}(N) ,$$

where  $N = \sum_{i=1}^k n_i$ , and `invalid-allocation` denotes the event in which `Allocation-Alg`( $n_1, \dots, n_k$ ) outputs an allocation such that there exists two different actual placements for some element, or that there exists some element with no actual placement. The probability is taken over the internal coin tosses of the algorithm `Allocation-Alg`.

- **Space:** For every input  $(n_1, \dots, n_k)$ , the array produced by `Allocation-Alg` is of size at most  $s(N)$ , where  $N = \sum_{i=1}^k n_i$ .
- **Locality:** For every input  $(N, n_i)$ , the algorithm `RangesGen` outputs at most  $d(N)$  ranges.
- **Read efficiency:** For every input  $(N, n_i)$  for the algorithm `RangesGen` it holds that:

$$\frac{\sum_{j=1}^d |b_j - a_j|}{n_i} \leq r(N) ,$$

where  $\{[a_1, b_1], \dots, [a_d, b_d]\} \leftarrow \text{RangesGen}(N, n_i)$ .

**Concrete instantiations.** Our two-dimensional balanced allocation schemes described in Section 3 have the above-discussed two-step structure (more accurately, they can be easily re-written as algorithms that follow this structure). From Corollaries 3.3 and 3.6 we obtain that our one-choice scheme `OneChoiceAllocation` (Algorithm 3.1) is an  $(O(N), O(1), \tilde{O}(\log N))$ -allocation algorithm, and that our two-choice scheme `TwoChoiceAllocation` (Algorithm 3.4) is an  $(O(N), O(1), \tilde{O}(\log \log N))$ -allocation algorithm for databases in which no list has more than  $N^{1-1/\log \log N}$  elements.

## 4.2 From Allocation Algorithms to SSE Schemes

We show a generic transformation from any allocation algorithm `Allocation-Alg` = (`RangesGen`, `Allocation`) to a searchable symmetric encryption scheme. In our SSE scheme, the client will run the `RangesGen` and the `Allocation` procedures, and then encrypt each of the document identifiers from list  $\text{DB}(w)$  with a key that is derived from the keyword  $w$  using a pseudorandom function. In addition, any unused entry in the array will be filled with a uniform random string of the appropriate length. Then, when issuing a query corresponding to a keyword  $w$ , the client will ask the server to retrieve the encrypted content of all possible locations of the list  $\text{DB}(w)$ . Since these locations are chosen independently at random, this does not reveal any information on the list except for its length. The client then decrypts the contents using the respective key for the list. The formal construction is described in Construction 4.3.

**Theorem 4.4.** *Assume that PRF is a pseudorandom function, and that (Enc, Dec) has pseudorandom ciphertexts, as well as elusive and efficiently verifiable range. Moreover, let Allocation-Alg be an  $(s, d, r)$ -allocation algorithm. Then, Construction 4.3 is an adaptive  $\mathcal{L}_{\text{size}}^{\text{adap}}$ -secure searchable symmetric encryption scheme, with space  $O(s(N))$ , locality  $O(d(N))$ , and read efficiency  $O(r(N))$ .*

**Proof.** The correctness of the scheme follows from the elusive and verifiable range properties of the underlying encryption scheme, and from the correctness of the procedure `Allocation-Alg`. In particular, from the correctness of the latter procedure and the pseudorandomness of PRF (i.e., the  $r_i$ 's are computationally indistinguishable from uniform values), we have all placeholders for all elements of all lists in the array `map`. Then, we use this array to create the array `Data`. From the verifiable range property, encryptions under different keys do not “collide”. In addition, since the encryption scheme has elusive range, the random elements that the client introduces (see Eq. (1)) do not collide with valid encryptions of any of the actual elements. Furthermore, with all but a negligible probability, encryptions under one key are not valid for other keys, and the probability that the same label  $\ell_i$  or the same key  $K_i$  appear more than once is negligible due to the pseudorandomness of PRF.

Regarding the efficiency, the size of the encrypted database is  $|\text{Data}| + |\text{HT}|$ , where  $|\text{Data}|$  is of size  $s(N)$  and `HT` contains exactly  $N$  elements of at most  $\log N = O(\log \lambda)$  bits each (and thus can be stored using  $O(N)$  machine words in the unit-cost RAM model). Therefore, the space usage is  $O(s(N) + N) = O(s(N))$  (note that  $s(N) \geq N$  always hold). Locality and read efficiency are straightforward.

We now prove the adaptive security of the scheme with respect to the leakage function  $\mathcal{L}_{\text{size}}^{\text{adap}}$ . Recall that in both the real and the ideal executions, the adversary  $\mathcal{A}$  first outputs `DB`. In the real execution, it then receives the encrypted database  $\text{EDB} \leftarrow \text{EDBSetup}(K, \text{DB})$ . In the ideal execution, the (interactive) simulator  $\mathcal{S}$  receives  $\mathcal{L}_{\text{size}}^{\text{adap}}(\text{DB}) = N$  and has to generate `EDB`. Then,  $\mathcal{A}$  adaptively issue queries  $\mathbf{w}$ . For each query  $w_i \in \mathbf{w}$ , in the real execution it receives the token  $\tau_i \leftarrow \text{TokGen}(K, w_i)$ , and in the ideal execution the simulator  $\mathcal{S}$  receives  $\mathcal{L}_{\text{size}}^{\text{adap}}(\text{DB}, \{w_j\}_{j < i}, w_i) =$

**CONSTRUCTION 4.3.**

**Parameters:** A keyword set  $W = \{w_1, \dots, w_{n_W}\}$ , and a database  $DB = \{DB(w_1), \dots, DB(w_{n_W})\}$ . We denote by  $N = \sum_{i=1}^{n_W} |DB(w_i)|$  the size of the database, and for each keyword  $w_i$  we denote by  $n_i = |DB(w_i)|$  the number of documents containing  $w_i$ .

**Key generator.** The algorithm  $\text{KeyGen}(1^\lambda)$  samples a PRF key  $K$ .

**Setup.** The algorithm  $\text{EDBSetup}(DB, K)$  proceeds as follows:

1. Initialize an empty set  $S$ .
2. For every keyword  $w_i \in W$ :
  - (a) Let  $DB(w_i) = \{\text{id}_1^{(i)}, \dots, \text{id}_{n_i}^{(i)}\}$ .
  - (b) Compute  $(\ell_i, k_i, r_i, K_i) = \text{PRF}_K(w_i)$ .
  - (c) Compute  $\widehat{n}_i = n_i \oplus k_i$ .
  - (d) Add the pair  $(\ell_i, \widehat{n}_i)$  to the set  $S$ .
  - (e) Compute  $R_i = \text{RangesGen}(N, n_i; r_i)$ .
3. Pad  $S$  until it contains exactly  $N$  pairs with random elements. Uniformly shuffle  $S$  and compute  $\text{HT} \leftarrow \text{HTSetup}(S)$ .
4. Compute  $\text{map} \leftarrow \text{Allocation}(\{(n_i, R_i)\}_{i=1}^{n_W})$ . If  $\text{map} = \perp$  then abort and output  $\perp$ . Otherwise, define the data block  $\text{Data}$  of size  $s(N)$  as follows, where for every  $1 \leq t \leq s(N)$ :

$$\text{Data}[t] = \begin{cases} \text{Enc}_{K_i}(\text{id}_j^{(i)}) & \text{if } \text{map}[t] = (i, j) \\ U_\ell & \text{otherwise} \end{cases}, \quad (1)$$

where  $U_\ell$  denote a uniformly and independently sampled  $\ell$ -bit string for each entry.

5. Output:  $\text{EDB} = (\text{Data}, \text{HT})$ .

**Token generator.** The algorithm  $\text{TokGen}(K, w_i)$  computes the derived keys  $(\ell_i, k_i, r_i, K_i) = \text{PRF}_K(w_i)$ . It outputs  $\tau_i = (\ell_i, k_i, r_i)$  as the public token, and  $\rho_i = K_i$  as the secret state for the algorithm  $\text{Resolve}$ .

**Search.** The algorithm  $\text{Search}(\tau_i, \text{EDB})$ , with  $\tau_i = (\ell_i, k_i, r_i)$  and  $\text{EDB} = (\text{Data}, \text{HT})$  proceeds as follows:

1. Obtain  $\widehat{n}_i \leftarrow \text{HTGet}(\text{HT}, \ell_i)$  and compute  $n_i = \widehat{n}_i \oplus k_i$ .
2. Run  $\text{RangesGen}(N, n_i; r_i)$ . Let  $P = \{[a_1, b_1], \dots, [a_d, b_d]\}$  be the resulting ranges.
3. Output  $\widetilde{D} = \bigcup_{[a_j, b_j] \in P} \text{Data}[a_j, \dots, b_j]$ .

**Resolve.** The algorithm  $\text{Resolve}(\rho_i, \widetilde{D})$  with  $\rho_i = K_i$ , decrypts each element in  $\widetilde{D}$  using the key  $K_i$ , and returns the set for which the decryption was successful.

$|DB(w_i)|$  and has to generate  $\tau_i$ . The adversary  $\mathcal{A}$  should not be able to distinguish between the real and ideal executions. The simulator  $\mathcal{S}$  works as follows.

- **Input:** Initially,  $\mathcal{L}_{\text{size}}^{\text{dap}}(DB) = N$ .

- **The simulator:**

- **Initialization phase.**

1.  $\mathcal{S}$  creates data block  $\text{Data}$  containing  $s(N)$  elements, each chosen uniformly and independently at random.
2. It initializes a set  $S$  with  $N$  random elements  $(\ell_{i'}, \widehat{n}_{i'})$ .
3. It uniformly shuffles  $S$  and computes  $\text{HT} \leftarrow \text{HTSetup}(S)$ .
4.  $\mathcal{S}$  then outputs:  $\text{EDB} = (\text{Data}, \text{HT})$ .

- **Query.** With each query that the experiment performs, the simulator receives the leakage

$|\text{DB}(w_i)|$ .

1.  $\mathcal{S}$  chooses a pair  $(\ell_{i'}, \widehat{n}_{i'})$  from  $S$  and removes the pair.
2. It computes  $k_i = |\text{DB}(w_i)| \oplus \widehat{n}_{i'}$ .
3. It uniformly samples  $r_i$ .
4.  $\mathcal{S}$  then outputs:  $\tau_i = (\ell_{i'}, k_i, r_i)$ .

We now claim that no adversary can distinguish whether  $(\text{EDB}, \tau = \{\tau_i\}_{w_i \in \mathbf{w}})$  was generated in the real execution or by the simulator. We show this by the following hybrid experiments:

- **Hyb<sub>0</sub>**. This is the real execution, where the adversary  $\mathcal{A}$  receives  $\text{EDB}$  as output of  $\text{EDBSetup}$ , and  $\tau_i = (\ell_i, k_i, r_i)$  with  $(\ell_i, k_i, r_i, K_i) = \text{PRF}_K(w_i)$  for every query  $w_i$ .
- **Hyb<sub>1</sub>**. In this experiment, we run the real world execution, where instead of using  $\text{PRF}_K$  (where  $K \leftarrow \text{KeyGen}(1^\lambda)$ ) in Step 2b of Construction 4.3, we use a truly random function. Note that as a result, all the keys  $(\ell_i, k_i, r_i, K_i)$  are uniform (for every  $w_i \in \mathbf{W}$ ). In addition, note that the key  $K$  as outputted by  $\text{KeyGen}$  is redundant.
- **Hyb<sub>2</sub>**. In this experiment, we change the value  $\widehat{n}_i$  (in Step 2c of Construction 4.3) to be a uniform value of the appropriate length. Then, when generating a token we compute  $k_i = \widehat{n}_i \oplus n_i$  instead of using the key  $k_i$  generated by the truly random function. Note that this makes the key  $k_i$  generated by the truly random function redundant.
- **Hyb<sub>3</sub>**. Here, for every  $w_i \in \mathbf{W}$ , we replace the encryption of each element in  $\text{DB}(w_i)$  (i.e., Eq. (1) of Construction 4.3) with an independent truly random value of the appropriate length. As a result, all elements of the data block  $\text{Data}$  are uniform and independent.
- **Hyb<sub>4</sub>**. The last hybrid is the ideal execution. Here we run the simulator  $\mathcal{S}$  defined above instead of the algorithms  $\text{EDBSetup}$  and  $\text{TokGen}$ .

$\text{Hyb}_0$  and  $\text{Hyb}_1$  are computationally indistinguishable from the security assumption of the pseudorandom function  $\text{PRF}$ . The experiments  $\text{Hyb}_1$  and  $\text{Hyb}_2$  are identically distributed since for every  $w_i \in \mathbf{W}$ ,  $\widehat{n}_i$  is uniformly distributed in both experiments, and for every query  $w_i \in \mathbf{w}$ ,  $k_i$  is determined by  $k_i = \widehat{n}_i \oplus n_i$ .  $\text{Hyb}_2$  and  $\text{Hyb}_3$  are computationally indistinguishable based on the pseudorandom ciphertexts property of the encryption scheme  $\text{Enc}$ . Finally,  $\text{Hyb}_3$  and  $\text{Hyb}_4$  are statistically close, where the only difference is the possibility of failure of the algorithm  $\text{Allocation}$ , which is negligible in  $N$ . This concludes the proof.  $\blacksquare$

**Conclusion.** By combining Theorem 4.4 with the allocation algorithms  $\text{OneChoiceAllocation}$  (Algorithm 3.1) and  $\text{TwoChoiceAllocation}$  (Algorithm 3.4) we obtain the following corollary:

**Corollary 4.5.** *Assuming the existence of one-way functions, there exist searchable symmetric encryption schemes with the following properties:*

1. Space  $O(N)$ , locality  $O(1)$ , and read efficiency  $\tilde{O}(\log N)$  without any assumptions on the structure of the database.
2. Space  $O(N)$ , locality  $O(1)$ , and read efficiency  $\tilde{O}(\log \log N)$  assuming that no keyword appears in more than  $N^{1 - \frac{1}{\log \log N}}$  documents.

### 4.3 Extensions

We discuss the following extensions of our generic transformation from allocation algorithms to searchable symmetric encryption schemes (Construction 4.3).

**Reducing one round of interaction.** When considering the SSE scheme that is based on our one-choice allocation algorithm `OneChoiceAllocation`, we can in fact reduce one round of interaction (as in [CT14]). Specifically, we revisit Construction 4.3 and provide some modifications. We then claim that under some additional assumptions on the underlying allocation algorithm (which are satisfied by our one-choice algorithm), the resulting construction is still a secure SSE scheme.

The modifications are as follows. Assume that while the client encrypts the lists, in addition it shuffles the elements in each bin. Then, we just combine the `Resolve` algorithm into the algorithm `Search`. That is, we also give the server the secret state  $\rho_i$ , so that it can perform the `Resolve` algorithm by itself and will not have to send the additional message to the client. Given the previous algorithms (`KeyGen`, `EDBSetup`, `TokGen`, `Search`, `Resolve`), our new scheme defines `EDBSetup` algorithm is modified as above (i.e., shuffles each bin), and for  $(\tau_i, \rho_i) = \text{TokGen}(K, w_i)$ , we define the new search algorithm as  $\text{Search}'(\tau_i, \text{EDB}) = \text{Resolve}(\rho_i, \text{Search}(\tau_i, \text{EDB}))$ . The scheme omits the `Resolve` algorithm, since the server already obtains the list of identifiers (or, documents).

We note that this is secure since the allocation algorithm satisfies the following (strong) property: not only that the *possible* location of each list is independent of the possible locations of other lists, but rather, the *actual* location of each list is independent of the structure of the other lists, conditioned that the allocation was successful (no bin was overflowing, which occurs with all but negligible probability). However, the security obtained in this manner is only static in the standard model (or adaptive in the random oracle model) – see Section 5 where we deal with a similar issue and provide a more elaborate discussion.

**On the actual leakage of the scheme.** The security obtained in Theorem 4.4 is for the modest leakage function  $\mathcal{L}_{\text{sizes}}^{\text{adap}}$ . However, this holds as long as the client stores a database of keyword-document pairs, and the decryption is done by the client and is not revealed to the server. In case that the database consists of keyword-identifier pairs, then recall that in the second round the client sends the identifiers and receives back the encrypted documents. As a result, although the first round of interaction can be simulated using the  $\mathcal{L}_{\text{sizes}}^{\text{adap}}$ -leakage function, this leakage is not enough for simulating the second round, and the actual leakage is  $\mathcal{L}_{\text{min}}^{\text{adap}}$ .

#### 4.4 A Limitation of Allocation Algorithms

We prove that any “very efficient” allocation algorithm that follows Structure 4.1, must assume a certain upper bound on the length of the longest list (this bound essentially matches the bound that is assumed by our two-choice algorithm). This is a generalization of Section 3.2.2, where here we consider some more general algorithms, but still follow essentially the same approach.

We note that since we consider a specific and rather restricted structure of allocation algorithms, we are able to prove a tighter bound than the one presented by Cash and Tessaro [CT14]. In fact, our construction in Section 5, which is based on a different technique (and not on allocation algorithms), circumvents this lower bound. This shows some dichotomy with respect to allocation algorithms: While these algorithms behave nicely while restricting the length of the maximal list (with respect to all three measures of efficiency), their efficiency becomes somewhat unsatisfying without this restriction. Adding the same restriction on the length of the maximal list to any of the other SSE schemes that have appeared in the literature (as well as to our construction in Section 5), does not seem to improve these constructions by much and they all behave asymptotically the same<sup>8</sup>.

---

<sup>8</sup>For example, assuming that there are no lists of length greater than  $N^{1-1/\log \log N}$ , the space usage of our construction in Section 5 is reduced from  $O(N \log N)$  to  $O(N \log N(1 - 1/\log \log N))$ , which is still  $O(N \log N)$ , while its locality and read efficiency remain exactly the same.

Our argument here is based on that from Section 3.2.2, where we generalize it in the following two directions. First, the `RangesGen` algorithm may choose an arbitrary number of ranges using some arbitrary distributions, and it is not restricted to choosing only two ranges of the appropriate length uniformly at random. For instance, the `RangesGen` may output  $d$  dependent ranges for the same list. Second, in our two-choice allocation scheme we analyzed some specific `Allocation` algorithm. Here, we show that the limitation holds for any implementation of the `Allocation` algorithm.

**Proposition 4.6.** *Let `Allocation-Alg` = (`RangesGen`, `Allocation`) be an  $(s, d, r)$ -allocation algorithm with  $s(N) = O(N \log N)$  and  $d(N)^2 r(N) \leq \log N / (c \cdot \log \log N)$  for some constant  $c \geq 1$ , and let  $\epsilon(N) = c^{-1} \cdot d(N)^{-2} r(N)^{-1}$ . Then, for every  $N$ , there exists an input  $(n_1, \dots, n_k)$  with  $N = \sum_{i=1}^k n_i$  and  $\max n_i = N^{1-\epsilon}$ , such that `Allocation-Alg`  $(n_1, \dots, n_k)$  fails with a noticeable probability.*

**Proof.** We modify `RangesGen` such that it returns exactly  $d$  ranges, each of size exactly  $r \cdot n_i$ . Surely this only increases the probability that we will be able to accommodate the lists. Similarly to the proof of Claim 3.9, we take  $L \leq N^\epsilon$  lists of size of  $n' = N^{1-\epsilon}$  for some  $\epsilon$ , and an arbitrary number of lists (and structure) for the remaining  $N - L \cdot N^{1-\epsilon}$  elements (but still with no list longer than  $n'$ ). From now on we will only refer to the first  $L$  lists.

For  $1 \leq \ell \leq L$ ,  $1 \leq j \leq d$  and  $1 \leq i \leq s$ , we define a random variable  $X_{\ell,j}^i \in \{0, 1\}$ , which indicates whether the  $j$ th range of the  $\ell$ th list contains the  $i$ th memory cell or not. We also define  $W_\ell = (X_{\ell,j}^i)_{j,i}$  of size  $d \times s$ , which gathers the variables that are related to the  $\ell$ th list. We point out that the variables  $W_1, \dots, W_L$  are independent and identically distributed (IID).

We now compute the probability that all ranges that are chosen to the  $L$  lists share the same ranges, or at least, share some few memory cells. As we will see, this will imply that there are a lot of intersections between the ranges and therefore for appropriate choice of parameters, we cannot accommodate the lists.

We denote the event  $\mathcal{E}_j^i = \{X_{\ell,j}^i = 1 : \forall 1 \leq \ell \leq L\}$ , i.e., the  $j$ th range of all lists contains the  $i$ th memory cell. For all  $\ell$  we have  $\sum_{i=1}^s X_{\ell,1}^i = r \cdot n'$ , so  $\sum_{i=1}^s \Pr[X_{\ell,1}^i = 1] = \sum_{i=1}^s E[X_{\ell,1}^i] = r \cdot n'$ , therefore there exists  $i_1$  with  $\Pr[X_{\ell,1}^{i_1} = 1] \geq \frac{r \cdot n'}{s}$ . This  $i_1$  is the same for all the lists because  $W_1, \dots, W_L$  are identically distributed. Thus, by independency (of choices between different lists) we get

$$\Pr[\mathcal{E}_1^{i_1}] = \prod_{\ell=1}^L \Pr[X_{\ell,1}^{i_1} = 1] \geq \left(\frac{r \cdot n'}{s}\right)^L = \left(\frac{N}{s} \cdot \frac{r}{N^\epsilon}\right)^L.$$

Let  $1 < j \leq d$  and suppose that we already defined  $i_1, \dots, i_{j-1}$ . Conditioned on  $\mathcal{E}_1^{i_1}, \dots, \mathcal{E}_{j-1}^{i_{j-1}}$ , the variables  $W_1, \dots, W_L$  are still IID, so applying the same argument, but with respect the conditional distribution, yields that there exists  $i_j$  such that

$$\Pr[\mathcal{E}_j^{i_j} \mid \mathcal{E}_1^{i_1}, \dots, \mathcal{E}_{j-1}^{i_{j-1}}] \geq \left(\frac{N}{s} \cdot \frac{r}{N^\epsilon}\right)^L.$$

Setting  $\mathcal{E} = \bigcap_{j=1}^d \mathcal{E}_j^{i_j}$  we get

$$\Pr[\mathcal{E}] = \prod_{j=1}^d \Pr[\mathcal{E}_j^{i_j} \mid \mathcal{E}_1^{i_1}, \dots, \mathcal{E}_{j-1}^{i_{j-1}}] \geq \left(\frac{N}{s} \cdot \frac{r}{N^\epsilon}\right)^{d \cdot L}.$$

Now if  $\mathcal{E}$  occurs then we have strictly less than  $2d \cdot r \cdot n'$  available memory cells to accommodate the lists. This is due to the fact that we can only use cells in the ranges  $[i_j - r \cdot n', i_j + r \cdot n' - 1]$  for  $1 \leq j \leq d$  (this is the point in the proof where we use the consecutiveness of the ranges). On

the other hand, the lists require  $L \cdot n'$  memory cells, so if we take  $L = 2d \cdot r$  we get that this event means failure for the Allocation algorithm. We conclude

$$\Pr[\text{Allocation-Alg}(n_1, \dots, n_k) \text{ fails}] \geq \Pr[\mathcal{E}] \geq \left(\frac{N}{s} \cdot \frac{r}{N^\epsilon}\right)^{2d^2r} = 2^{-2d^2r(\epsilon \log N + \log \frac{s}{N} - \log r)} \quad (2)$$

so if  $d^2r(\epsilon \log N + \log \frac{s}{N} - \log r) = O(\log N)$  then this algorithm fails with a noticeable probability. Note that we need to have  $2d \cdot r \leq N^\epsilon$ , since otherwise by taking  $L = 2d \cdot r$  lists we will exceed  $N$  elements.

Finally, we show that under our assumptions on  $s, d, r, \epsilon$ , the probability of failure of Eq. (2) is noticeable and that  $2d \cdot r \leq N^\epsilon$ . We start with the failure probability. We have  $\epsilon = c \cdot d^{-2}r^{-1}$ ,  $d^2r \leq \log N / (c \cdot \log \log N)$  and  $\frac{s}{N} = O(\log N)$ . So we get that

$$d^2r \cdot \left(\epsilon \log N + \log \frac{s}{N} - \log r\right) \leq c^{-1} \cdot \log N + d^2r \cdot \log \frac{s}{N} = O(\log N).$$

It remains to show that  $2d \cdot r \leq N^\epsilon$  for a large enough  $N$ . This is equivalent to showing that  $\epsilon^{-1}(1 + \log d + \log r) \leq \log N$ . Indeed, we have that

$$\begin{aligned} \epsilon^{-1} \cdot (1 + \log d + \log r) &= d^2r \cdot (1 + \log d + \log r) \\ &\leq c \cdot d^2r \cdot (1 + \log(d^2r)) \\ &\leq \frac{\log N}{\log \log N} \cdot \left(1 + \log \left(\frac{\log N}{\log \log N}\right)\right) \\ &= \log N - \left(1 - \frac{1}{\log \log \log N}\right) \cdot \frac{\log N \log \log \log N}{\log \log N} \\ &\leq \log N, \end{aligned}$$

where the last inequality holds for  $N \geq 16$ . ■

**Remark.** Suppose that every range returned by  $\text{RangesGen}(N, n_i)$  has size at most  $f(N) \cdot n_i$  (clearly we can always take  $f(N) \leq r(N)$ ). It can be easily seen from the proof, that we can strengthen Proposition 4.6 by replacing each appearance of  $r(N)$  by  $f(N)$ , that is, we require  $d^2(N)f(N) \leq \log N / (c \cdot \log \log N)$  and set  $\epsilon(N) = c^{-1} \cdot d(N)^{-2}f(N)^{-1}$ .

## 5 Improving the Cash-Tessaro Scheme

In this section we improve the scheme proposed by Cash and Tessaro [CT14]. Their scheme offers space  $O(N \log N)$ , locality  $O(\log N)$ , and read efficiency  $O(1)$ . Our improvement reduces the locality from  $O(\log N)$  to  $O(1)$  while not hurting the space overhead or read efficiency. We begin by first describing the Cash-Tessaro scheme, and then explain our modification.

**The Cash-Tessaro scheme.** Let  $\text{DB} = \{\text{DB}(w_1), \dots, \text{DB}(w_{n_W})\}$  be a database of size  $N = \sum_{i=1}^{n_W} |\text{DB}(w_i)|$ , where each keyword  $w_i$  appears in  $n_i = |\text{DB}(w_i)|$  documents. The scheme consists of  $\log N$  hash tables, each of size  $O(N)$ , where the  $k$ th table will store  $N/2^k$  blocks of length  $2^k$ .

For every keyword  $w_i$ , let  $(s_{\log N}^{(i)}, \dots, s_0^{(i)})$  be the binary representation of  $n_i$ . Then, each list  $\text{DB}(w_i)$  is broken into at most  $\log N$  consecutive blocks, where the  $k$ th block consists of  $2^k$  document identifiers for every  $k$  for which  $s_k^{(i)} = 1$ . For every  $k$  we then insert all blocks of length  $2^k$  into the  $k$ th hash table. Note that there can be at most  $N/2^k$  such blocks, and each hash table is padded

with “dummy” blocks to contain exactly  $N/2^k$  blocks. The blocks corresponding to each list  $\text{DB}(w_i)$  are then encrypted using a key that is derived from  $w_i$ .

Reconstruction of a list  $\text{DB}(w_i)$  can be done easily using the key that is derived from  $w_i$  by querying each of the  $\log N$  hash tables. Therefore, the scheme uses space  $O(N \log N)$ , locality  $O(\log N)$  (since one has to query all  $\log N$  hash tables), and read efficiency  $O(1)$ .

**Our scheme.** As in the Cash-Tessaro scheme, our scheme consists of  $\log N$  hash tables of size  $N$  each, where the  $k$ th tables stores  $2^k$  blocks of length  $N/2^k$ . The main idea underlying our scheme is to store each list  $\text{DB}(w_i)$  in a single table, instead of splitting it across the  $\log N$  tables. Specifically, every list  $\text{DB}(w_i)$  of length  $n_i$  is stored in the  $p_i$ th table, where  $2^{p_i-1} < |\text{DB}(w_i)| \leq 2^{p_i}$ . We pad the list with dummy elements until its length is exactly  $2^{p_i}$ , and this increases the overall space by a factor of 2, but the entire list is now stored in one continuous block. In addition, we include a separate hash table that contains an encryption of  $p_i$  for every keyword  $w_i$ . We therefore obtain a scheme with space  $O(N \log N)$ , locality  $O(1)$ , and read efficiency  $O(1)$ . The scheme is presented as Construction 5.1.

**CONSTRUCTION 5.1.**

**Parameters.** Let  $t$  denote the value for which  $2^{t-1} < |\text{DB}| = N \leq 2^t$ . For every word  $w_i \in \mathcal{W}$ , let  $\text{DB}(w_i) = \{\text{id}_1, \dots, \text{id}_{n_i}\}$ , and let  $p_i$  denote the value for which  $2^{p_i-1} < n_i \leq 2^{p_i}$ .

**Key generator.** Algorithm  $\text{KeyGen}(1^\lambda)$  samples a key  $K$  for the PRF.

**Setup.** Algorithm  $\text{EDBSetup}(\text{DB}, K)$  proceeds as follows:

1. Initialize  $t + 1$  empty sets  $T_0, \dots, T_t$ , where  $T_j$  will store at most  $2^{t-j}$  pairs  $(\ell, d)$ , where  $|d| = 2^j$ . Initialize a set  $S$ , which will be converted to a lookup table.
2. For every keyword  $w_i \in \mathcal{W}$

(Let  $\text{DB}(w_i) = \{\text{id}_1, \dots, \text{id}_{n_i}\}$ , and recall that  $2^{p_i-1} < n_i \leq 2^{p_i}$ . If necessary, pad  $\text{DB}(w_i)$  with dummy identifiers in order to contain exactly  $2^{p_i}$  elements.)

- (a) Compute  $\text{PRF}_K(w_i) = ((\ell_i, K_i), (\widehat{\ell}_i, \widehat{K}_i))$ .
- (b) Create the blocks:

$$d_i = (\text{Enc}_{K_i}(\text{id}_1), \dots, \text{Enc}_{K_i}(\text{id}_{2^{p_i}})) \quad \text{and} \quad \widehat{n}_i = \text{Enc}_{\widehat{K}_i}(n_i).$$

- (c) Insert the pair  $(\ell_i, d_i)$  to the set  $T_{p_i}$  and  $(\widehat{\ell}_i, \widehat{n}_i)$  to the set  $S$ .

3. Pad the list  $S$  to contain exactly  $N$  elements by adding random elements, and pad each set  $T_j$  to contain exactly  $2^{t-j}$  pairs by adding random elements.
4. For each set  $S, T_0, \dots, T_t$ , uniformly shuffle the set, generate the respective hash table by invoking the  $\text{HTSetup}$  algorithm, and obtain the hash tables  $\text{HT}(S), (\text{HT}(L_0), \dots, \text{HT}(L_t))$ .
5. Output  $\text{EDB} = (\text{HT}(S), (\text{HT}(L_0), \dots, \text{HT}(L_t)))$ .

**Token generator.** Algorithm  $\text{TokGen}(K, w_i)$  outputs the four derived keys

$$\tau_i = ((\ell_i, K_i), (\widehat{\ell}_i, \widehat{K}_i)) = \text{PRF}_K(w_i).$$

**Search.** Algorithm  $\text{Search}(\tau_i, \text{EDB})$ , where  $\tau_i = ((\ell_i, K_i), (\widehat{\ell}_i, \widehat{K}_i))$  and  $\text{EDB} = (\text{HT}(S), \text{HT}(L_0), \dots, \text{HT}(L_t))$  works as follows:

1. Invoke  $\text{HTLookup}$  on the hash table  $\text{HT}(S)$  and label  $\widehat{\ell}_i$ , receive back  $\widehat{n}_i = \text{Enc}_{\widehat{K}_i}(n_i)$ . Decrypt  $n_i$  using the key  $\widehat{K}_i$ , and compute  $p_i$ .
2. Invoke  $\text{HTLookup}$  on the hash table  $\text{HT}(T_{p_i})$  and the label  $\ell_i$ . Obtain the block  $d_i = (\text{Enc}_{K_i}(\text{id}_1), \dots, \text{Enc}_{K_i}(\text{id}_{2^{p_i}}))$ . Decrypt the first  $n_i$  elements of this block using the key  $K_i$  and return them.

**Theorem 5.2.** Assume that PRF is a pseudorandom function, and that Enc has pseudorandom



ciphertexts. Then, Construction 5.1 is a static  $\mathcal{L}_{\min}$ -secure searchable symmetric encryption scheme for databases of size  $N$  with space  $O(N \log N)$ , locality  $O(1)$ , and read efficiency  $O(1)$ .

**Proof.** It is easy to see that the scheme is correct except for the event where the same label appears more than once, which occurs in some negligible probability. For proving the security of the scheme, recall that in the real execution, the adversary  $\mathcal{A}$  outputs  $\text{DB}, \mathbf{w}$ , and is given the encrypted database  $\text{EDB}$  and the tokens  $\tau_i = \text{TokGen}(K, w_i)$  for every  $w_i \in \mathbf{w}$ . In the ideal execution, the simulator  $\mathcal{S}$  is given only  $\mathcal{L}_{\min}(\text{EDB}, \mathbf{w})$  and should output both  $\text{EDB}$  and  $\boldsymbol{\tau} = \{\tau_i\}_{w_i \in \mathbf{w}}$  such that the adversary cannot distinguish between these two executions. Observe that the algorithms  $\text{KeyGen}, \text{EDBSetup}$  and  $\text{TokGen}$  are not invoked in the ideal execution (by the experiment itself. They may be invoked by the simulator).

The simulator  $\mathcal{S}$  works as follows:

- **Input:**  $\mathcal{L}(\text{DB}, \mathbf{w}) = (N, \{\text{DB}(w_i)\}_{w_i \in \mathbf{w}})$ .
- **The simulator:**
  1. The simulator  $\mathcal{S}$  chooses (uniformly) random keys  $\tau_i = ((\ell_i, K_i), (\widehat{\ell}_i, \widehat{K}_i))$  for every  $w_i \in \mathbf{w}$ .
  2.  $\mathcal{S}$  initializes the sets  $L_0, \dots, L_t$  and  $S$  to be empty sets. For every  $w_i \in \mathbf{w}$  and given  $\text{DB}(w_i)$ , it computes  $n_i$  and  $p_i$ . If necessary, it pads  $\text{DB}(w_i)$  to size  $2^{p_i}$  with some random additional id's.
  3. For every  $w_i \in \mathbf{w}$ ,  $\mathcal{S}$  computes:

$$d_i = (\text{Enc}_{K_i}(\text{id}_1), \dots, \text{Enc}_{K_i}(\text{id}_{2^{p_i}})) \quad \text{and} \quad \widehat{n}_i = \text{Enc}_{\widehat{K}_i}(n_i).$$

It adds the pair  $(\widehat{\ell}_i, \widehat{n}_i)$  to the set  $S$  and  $(\ell_i, d_i)$  to the set  $L_{p_i}$ .

4. After adding all the keywords in  $\mathbf{w}$ ,  $\mathcal{S}$  adds for each set  $L_j$  random elements until it contains exactly  $2^{t-j}$  elements, and it adds to  $S$  random elements until it contains exactly  $N$  elements.
5.  $\mathcal{S}$  uniformly shuffles each one of the sets  $S, T_0, \dots, T_t$ . It then generates hash table from each set using the algorithm  $\text{HTSetup}$ , and defines  $\text{EDB} = (\text{HT}(S), (\text{HT}(T_0), \dots, \text{HT}(T_t)))$ .
6.  $\mathcal{S}$  outputs  $\text{EDB}$  and  $\boldsymbol{\tau} = \{\tau_i\}_{w_i \in \mathbf{w}}$ .

We now claim that the adversary cannot distinguish between the pair  $(\text{EDB}, \boldsymbol{\tau})$  that it receives in the real execution and in the ideal execution. Consider the following hybrid experiments:

- **Hyb<sub>0</sub>.** This is the real execution, where the adversary  $\mathcal{A}$  receives  $\text{EDB}$  and  $\tau_i = \text{PRF}_K(w_i)$  for every  $w_i \in \mathbf{w}$ .
- **Hyb<sub>1</sub>.** This experiment is obtained from  $\text{Hyb}_0$  by replacing the pseudorandom function  $\text{PRF}_K(\cdot)$  with a truly random function. Observe that, as a result, the key  $K$  that is produced by  $\text{KeyGen}$  is redundant and the elements  $\tau_j = ((\ell_j, K_j), (\widehat{\ell}_j, \widehat{K}_j))$  are distributed uniformly.
- **Hyb<sub>2</sub>.** This experiment is obtained from  $\text{Hyb}_1$  as follows: For every  $w_j \in W \setminus \mathbf{w}$  we replace the values  $d_j$  and  $\widehat{n}_j$  (in Step 2b of  $\text{EDBSetup}$ ) with independent and uniformly-distributed values of the appropriate length.
- **Hyb<sub>3</sub>.** This is the ideal execution, where we run the simulator  $\mathcal{S}$  defined above.

We observe that  $\text{Hyb}_0$  and  $\text{Hyb}_1$  are computationally indistinguishable based on the security of the pseudorandom function  $\text{PRF}$ , and that  $\text{Hyb}_1$  and  $\text{Hyb}_2$  are computationally indistinguishable based on the pseudorandom ciphertexts property of the encryption scheme  $\text{Enc}$ . Finally, the experiments  $\text{Hyb}_2$  and  $\text{Hyb}_3$  are identical by the definition of the simulator  $\mathcal{S}$ . ■

**Adaptive security in the random-oracle model.** The scheme can be proven with adaptive security in the programmable random oracle model, using similar techniques as in [CT14]. Towards this goal, all we need is to instantiate the encryption scheme  $(\text{Enc}, \text{Dec})$  with an encryption that allows “explaining” a random ciphertext as the encryption of any particular message (as in [CT14]). Concretely, we can use  $\text{Enc}_K(x; r) = (r, H(K||r) \oplus x0^\lambda)$ , where  $H$  is the random oracle. Then, for any derived key  $k$ , random  $(r, c)$ , and a message  $m$ , we can program the random oracle such that  $\text{Enc}_k(m) = (r, c)$ , by fixing  $H(k||r) = c \oplus m0^\lambda$ .

For the proof of security, the simulator in the first step creates random sets  $L_0, \dots, L_t$  and  $S$ , apply the HTSetup on each one of them and outputs the resulting tables as EDB. Later, upon a query  $w_i$  with a list  $\text{DB}(w_i) = \{\text{id}_1, \dots, \text{id}_{n_i}\}$ , it computes  $p_i$ , chooses some random pair  $(\ell_i, (d_i^{(1)}, \dots, d_i^{(2^{p_i})}))$  from the list  $L_{p_i}$  (and removes this pair from  $L_{p_i}$ ), and chooses some random pair  $(\widehat{\ell}_i, \widehat{n}_i)$  from the set  $S$  (and removes it from  $S$ ). Then, it chooses random keys  $K_i, \widehat{K}_i$ , and programs the random oracle such that the following hold:

$$\text{Enc}_{K_i}(\text{id}_j) = d_i^{(j)} \quad \forall 1 \leq j \leq n_i, \quad \text{and} \quad \text{Enc}_{\widehat{K}_i}(n_i) = \widehat{n}_i.$$

If the programming cannot succeed (because the corresponding values are already set for  $H$ ), the simulator aborts. Otherwise, it outputs  $\tau_i = ((\ell_i, K_i), (\widehat{\ell}_i, \widehat{K}_i))$ . We omit a formal analysis of the above, as it follows standard techniques. We conclude the following theorem:

**Theorem 5.3.** *Construction 5.1 is adaptive  $\mathcal{L}_{\min}^{\text{adap}}$ -secure searchable encryption scheme for database of size  $N$ , with space  $O(N \log N)$ , locality  $O(1)$  and read efficiency  $O(1)$  in the random-oracle model.*

**Adaptive security in the plain model.** We note that the usage of random oracle can be avoided using somewhat similar ideas to our solution in Section 4.2. That is, we modify the scheme such that the token that the client sends does not give the server the ability to decrypt the elements of the resulting list, and these decryptions are performed by the client (i.e., the client does not send the key  $K_i$ ). The encryption of  $n_i$  is performed using a one-time pad, similarly to Construction 4.3. However, in case where the client stores database of keywords/identifiers (and not keywords/documents), this results in one additional round of interaction, where in the second round the client sends the decrypted list as a list of documents’ identifiers the server has to fetch from the database.

## Acknowledgments

We thank Raphael Bost, Angèle Bossuat, Pierre-Alain Fouque and Brice Minaud for observing that our bound on the probability of the event  $\neg \mathcal{E}_0$  in the proof of Theorem 3.5 was incorrect in a previous version of this paper.

## References

- [ABK<sup>+</sup>99] Y. Azar, A. Z. Broder, A. R. Karlin, and E. Upfal. Balanced allocations. *SIAM Journal on Computing*, 29(1):180–200, 1999.
- [ANS10] Y. Arbitman, M. Naor, and G. Segev. Backyard cuckoo hashing: Constant worst-case operations with a succinct representation. In *Proceedings of the 51st Annual IEEE Symposium on Foundations of Computer Science*, pages 787–796, 2010.

- [BM05] A. Z. Broder and M. Mitzenmacher. Multidimensional balanced allocations. In *Proceedings of the 15th Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 195–196, 2005.
- [CGK<sup>+</sup>06] 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*, pages 79–88, 2006.
- [CGK<sup>+</sup>11] R. Curtmola, J. A. Garay, S. Kamara, and R. Ostrovsky. Searchable symmetric encryption: Improved definitions and efficient constructions. *Journal of Computer Security*, 19(5):895–934, 2011.
- [CGP<sup>+</sup>15] D. Cash, P. Grubbs, J. Perry, and T. Ristenpart. Leakage-abuse attacks against searchable encryption. In *Proceedings of the 22nd ACM Conference on Computer and Communications Security*, pages 668–679, 2015.
- [CJJ<sup>+</sup>13] 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 '13*, pages 353–373, 2013.
- [CJJ<sup>+</sup>14] D. Cash, J. Jaeger, S. Jarecki, C. S. Jutla, H. Krawczyk, M. Rosu, and M. Steiner. Dynamic searchable encryption in very-large databases: Data structures and implementation. In *Proceedings of the 21st Annual Network and Distributed System Security Symposium*, 2014.
- [CK10] M. Chase and S. Kamara. Structured encryption and controlled disclosure. In *Advances in Cryptology - ASIACRYPT '10*, pages 577–594, 2010.
- [CM05] Y.-C. Chang and M. Mitzenmacher. Privacy preserving keyword searches on remote encrypted data. In *Proceedings of the 3rd International Conference on Applied Cryptography and Network Security*, pages 442–455, 2005.
- [CT14] D. Cash and S. Tessaro. The locality of searchable symmetric encryption. In *Advances in Cryptology - EUROCRYPT '14*, pages 351–368, 2014.
- [DP08] M. Dietzfelbinger and R. Pagh. Succinct data structures for retrieval and approximate membership. In *Proceedings of the 35th International Colloquium on Automata, Languages and Programming*, pages 385–396, 2008.
- [DP09] D. P. Dubhashi and A. Panconesi. Concentration of Measure for the Analysis of Randomized Algorithms. Cambridge University Press, 2009.
- [DR98] D. Dubhashi and D. Ranjan. Balls and bins: A study in negative dependence. *Random Structures & Algorithms*, v.13, 99-124 (1998), 13:99–124, 09 1998.
- [Goh03] E. Goh. Secure indexes. Cryptology ePrint Archive, Report 2003/216, 2003.
- [Hag98] T. Hagerup. Sorting and searching on the word RAM. In *Proceedings of the 15th Annual Symposium on Theoretical Aspects of Computer Science*, pages 366–398, 1998.
- [HMP01] T. Hagerup, P. B. Miltersen, and R. Pagh. Deterministic dictionaries. *Journal of Algorithms*, 41(1):69–85, 2001.

- [JDP83] K. Joag-Dev and F. Proschan. Negative association of random variables with applications. *The Annals of Statistics*, 11(1):286–295, 1983.
- [KO12] K. Kurosawa and Y. Ohtaki. UC-secure searchable symmetric encryption. In *Proceedings of the 16th International Conference on Financial Cryptography and Data Security*, pages 285–298, 2012.
- [KO13] K. Kurosawa and Y. Ohtaki. How to update documents verifiably in searchable symmetric encryption. In *Proceedings of the 12th International Conference on Cryptology and Network Security*, pages 309–328, 2013.
- [KP13] S. Kamara and C. Papamanthou. Parallel and dynamic searchable symmetric encryption. In *Proceedings of the 16th International Conference on Financial Cryptography and Data Security*, pages 258–274, 2013.
- [KPR12] S. Kamara, C. Papamanthou, and T. Roeder. Dynamic searchable symmetric encryption. In *Proceedings of the 19th ACM Conference on Computer and Communications Security*, pages 965–976, 2012.
- [KS81] A. Khursheed and K. Saxena. Positive dependence in multivariate distributions. *Communications in Statistics - Theory and Methods*, 10:1183–1196, 01 1981.
- [LP09] Y. Lindell and B. Pinkas. A proof of security of Yao’s protocol for two-party computation. *Journal of Cryptology*, 22(2):161–188, 2009.
- [Mil99] P. B. Miltersen. Cell probe complexity - a survey. In *Proceedings of the 19th Conference on the Foundations of Software Technology and Theoretical Computer Science, Advances in Data Structures Workshop*, 1999.
- [MRS00] M. Mitzenmacher, A. W. Richa, and R. Sitaraman. The power of two random choices: A survey of techniques and results. In *Handbook of Randomized Computing*, pages 255–312, 2000.
- [MU05] M. Mitzenmacher and E. Upfal. Probability and computing – randomized algorithms and probabilistic analysis. Cambridge University Press, 2005.
- [PP08] A. Pagh and R. Pagh. Uniform hashing in constant time and optimal space. *SIAM Journal on Computing*, 38(1):85–96, 2008.
- [PR04] R. Pagh and F. F. Rodler. Cuckoo hashing. *Journal of Algorithms*, 51(2):122–144, 2004.
- [SWP00] D. X. Song, D. Wagner, and A. Perrig. Practical techniques for searches on encrypted data. In *Proceedings of the 21st Annual IEEE Symposium on Security and Privacy*, pages 44–55, 2000.
- [vLSD<sup>+</sup>10] P. van Liesdonk, S. Sedghi, J. Doumen, P. H. Hartel, and W. Jonker. Computationally efficient searchable symmetric encryption. In *Proceedings of 7th VLDB Workshop on Secure Data Management*, pages 87–100, 2010.
- [Waj17] D. Wajc. Negative association-definition , properties , and applications. Available at <https://www.cs.cmu.edu/~dwajc/notes/Negative%20Association.pdf>, 2017.