

# Tight Tradeoffs in Searchable Symmetric Encryption

Gilad Asharov\*

Gil Segev†

Ido Shahaf†‡

## Abstract

A searchable symmetric encryption (SSE) scheme enables a client to store data on an untrusted server while supporting keyword searches in a secure manner. Recent experiments have indicated that the practical relevance of such schemes heavily relies on the tradeoff between their *space overhead*, *locality* (the number of non-contiguous memory locations that the server accesses with each query), and *read efficiency* (the ratio between the number of bits the server reads with each query and the actual size of the answer). These experiments motivated Cash and Tessaro (EUROCRYPT '14) and Asharov et al. (STOC '16) to construct SSE schemes offering various such tradeoffs, and to prove lower bounds for natural SSE frameworks. Unfortunately, the best-possible tradeoff has not been identified, and there are substantial gaps between the existing schemes and lower bounds, indicating that a better understanding of SSE is needed.

We establish tight bounds on the tradeoff between the space overhead, locality and read efficiency of SSE schemes within two general frameworks that capture the memory access pattern underlying all existing schemes. First, we introduce the “pad-and-split” framework, refining that of Cash and Tessaro while still capturing the same existing schemes. Within our framework we significantly strengthen their lower bound, proving that any scheme with locality  $L$  must use space  $\Omega(N \log N / \log L)$  for databases of size  $N$ . This is a tight lower bound, matching the tradeoff provided by the scheme of Demertzis and Papamanthou (SIGMOD '17) which is captured by our pad-and-split framework.

Then, within the “statistical-independence” framework of Asharov et al. we show that their lower bound is essentially tight: We construct a scheme whose tradeoff matches their lower bound within an additive  $O(\log \log \log N)$  factor in its read efficiency, once again improving upon the existing schemes. Our scheme offers optimal space and locality, and nearly-optimal read efficiency that depends on the frequency of the queried keywords: For a keyword that is associated with  $n = N^{1-\epsilon(n)}$  document identifiers, the read efficiency is  $\omega(1) \cdot \epsilon(n)^{-1} + O(\log \log \log N)$  when retrieving its identifiers (where the  $\omega(1)$  term may be arbitrarily small, and  $\omega(1) \cdot \epsilon(n)^{-1}$  is the lower bound proved by Asharov et al.). In particular, for any keyword that is associated with at most  $N^{1-1/o(\log \log \log N)}$  document identifiers (i.e., for any keyword that is not exceptionally common), we provide read efficiency  $O(\log \log \log N)$  when retrieving its identifiers.

---

A preliminary version of this work appeared in *Advances in Cryptology - CRYPTO '18*.

\*Department of Computer Science, Bar-Ilan University, Israel. Email: [Gilad.Asharov@biu.ac.il](mailto:Gilad.Asharov@biu.ac.il). Most of the work was done while at Cornell Tech, supported by a Junior Fellow award from the Simons Foundation.

†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 European Union's Horizon 2020 Framework Program (H2020) via an ERC Grant (Grant No. 714253), 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 US-Israel Binational Science Foundation (Grant No. 2014632).

‡Supported by the Clore Israel Foundation via the Clore Scholars Programme.

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Our Contributions . . . . .	2
1.2	Overview of Our Contributions . . . . .	3
1.3	Related Work . . . . .	6
1.4	Paper Organization . . . . .	7
<b>2</b>	<b>Preliminaries</b>	<b>7</b>
2.1	Searchable Symmetric Encryption . . . . .	8
2.2	Private-Key Encryption . . . . .	11
2.3	Static Hash Tables . . . . .	11
2.4	Cuckoo Hashing with a Stash . . . . .	11
<b>3</b>	<b>The Pad-and-Split Framework: A Stronger Lower Bound</b>	<b>12</b>
3.1	The Pad-and-Split Framework . . . . .	12
3.2	The Generality of the Pad-and-Split Framework . . . . .	19
3.3	An Optimal Instantiation for any Locality . . . . .	20
3.4	Our Lower Bound for Pad-and-Split Schemes . . . . .	21
<b>4</b>	<b>The Statistical-Independence Framework: A Leveled Two-Choice Scheme</b>	<b>25</b>
4.1	The Statistical-Independence Framework . . . . .	26
4.2	Our Leveled Two-Choice Scheme . . . . .	28
	<b>References</b>	<b>30</b>

## 1 Introduction

A searchable symmetric encryption (SSE) scheme [SWP00, CGK<sup>+</sup>06] enables 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 data items that are associated with  $w$  (e.g., all document identifiers that contain  $w$ ). This typically consists of a two-stage process: First, the client encrypts her database and uploads it to the server, and then the client repeatedly queries the server with various keywords by providing the server with keyword-specific search tokens. Informally, the security requirement of SSE schemes asks that the server does not learn any information about keywords for which the client did not issue any queries.

**The practical relevance of SSE schemes.** Motivated by the increasingly-growing technological interest in outsourcing data to remote (and thus potentially untrusted) servers, a very fruitful line of research in the cryptography community focused on the design of SSE schemes (e.g., [SWP00, 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, ANS<sup>+</sup>16, DP17]). Most of the proposed schemes offer strong and meaningful notions of security, and some even extend the basic keyword search functionality to more expressive ones.

Despite these promising developments, Cash et al. [CJJ<sup>+</sup>13] showed via experiments with real-world databases that the practical performance of the known schemes is quite disappointing, and scales badly to large databases. Somewhat surprisingly, they observed that performance issues resulting from impractical memory layouts may be significantly more crucial compared to performance issues resulting from the cryptographic processing of the data. More specifically, Cash et al. observed that schemes with poor *locality* (i.e., schemes in which the server has to access a rather large number of *non-contiguous* memory locations with each query) have poor practical performance when dealing with large databases that require the usage of disk-storage mechanisms.

Practical locality, however, is obviously insufficient: Any practically-relevant SSE scheme should (at least) not suffer from either a significant *space overhead* (i.e., encrypted databases should not be much larger than the original databases), or from a poor *read efficiency* (i.e., servers should not read much more data than needed for answering each query)<sup>1</sup>.

**Efficiency tradeoffs and existing lower bounds.** This state of affairs naturally poses the challenge of constructing an SSE scheme that simultaneously enjoys asymptotically-optimal space overhead, locality, and read efficiency – but unfortunately no such scheme is currently known. This has motivated Cash and Tessaro [CT14] to initiate the study of understanding the tradeoff between these central measures of efficiency. They proved a lower bound showing that, for a large and natural class of SSE schemes, it is in fact impossible to simultaneously enjoy asymptotically-optimal space overhead, locality, and read efficiency. Specifically, they considered the class of SSE schemes with “non-overlapping reads”: Schemes in which distinct keywords induce non-overlapping memory regions which the server may access upon their respective queries (we refer the reader to the work of Cash and Tessaro [CT14] for a formal definition of their notion of non-overlapping reads).

The class of SSE schemes with non-overlapping reads captures the basic techniques underlying all existing SSE schemes other than two schemes proposed by Asharov et al. [ANS<sup>+</sup>16]. These two schemes may have arbitrary overlapping reads, and offer an improved tradeoff between their space overhead, locality, and read efficiency compared to the previously suggested schemes. This

---

<sup>1</sup>We consider the notions of locality and read efficiency as formalized by Cash and Tessaro [CT14]: The locality of a scheme is the 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.

tradeoff, however, is still non-optimal, and Asharov et al. showed that this is in fact inherent to their approach. Similarly to Cash and Tessaro, they proved that also for a different class of SSE schemes, it is impossible to simultaneously enjoy asymptotically-optimal space overhead, locality, and read efficiency. Specifically, they considered the class of SSE scheme with “statistically-independent reads”: Schemes in which distinct keywords induce statistically-independent memory regions which the server accesses upon their respective queries.

The lower bounds proved by Cash and Tessaro and by Asharov et al. capture all of the existing SSE schemes (except for various schemes with non-standard leakage or functionality that we do not consider in this work). That is, the basic techniques underlying each of the known SSE schemes belong either to the class of “non-overlapping reads” or to the class of “statistically-independent reads”. In both cases, however, the existing lower bounds are not tight, as there are still noticeable gaps between the lower bounds and the performance guarantees of the existing schemes (as we detail in the next section). This unsatisfying situation calls for obtaining a better understanding of SSE techniques: Either by strengthening the known lower bounds, or by designing new schemes with better performance guarantees.

## 1.1 Our Contributions

We prove tight bounds on the tradeoff between the space overhead, locality, and read efficiency of SSE schemes within the following two general frameworks:

**The pad-and-split framework:** We formalize a framework that refines the non-overlapping reads framework of Cash and Tessaro [CT14] while still capturing the same existing SSE schemes (i.e., all existing schemes other than those of Asharov et al. [ANS<sup>+</sup>16])<sup>2</sup>. We refer to this framework as the “pad-and-split” framework given the structure of the SSE schemes that it captures.

Within this framework we significantly strengthen the lower bound of Cash and Tessaro: We show that any pad-and-split scheme with locality  $L$  must use space  $\Omega(N \cdot \log N / \log L)$  for databases of size  $N$ . For example, for any constant locality (i.e.,  $L = O(1)$ ) and for any logarithmic locality (i.e.,  $L = O(\log N)$ ) our lower bound shows that any such scheme must use space  $\Omega(N \log N)$  and  $\Omega(N \log N / \log \log N)$ , respectively, and is thus not likely to be of substantial practical relevance (whereas the lower bound of Cash and Tessaro would only yield space  $\omega(N)$  when the locality is constant).

Then, we observe that our lower bound is in fact tight, as it is matched by a recent scheme proposed by Demertzis and Papamanthou [DP17] that is captured by our framework (i.e., their scheme is an optimal instantiation of our framework). We refer the reader to Sections 1.2 and 3 for a high-level overview and for a detailed description of this framework, its instantiations, and of our lower bound, respectively.

**The statistical-independence framework:** We consider the statistical-independence framework of Asharov et al. [ANS<sup>+</sup>16], and show that their lower bound for SSE schemes in this framework is essentially tight: Based on the existence of any one-way function, we construct a scheme whose efficiency guarantees match their lower bound for constant locality within an additive  $O(\log \log \log N)$  factor in the read efficiency, and improve upon those of their two schemes.

---

<sup>2</sup>Each of the schemes that are captured by our framework offers other important implementation details, improvements and optimizations that we do not intend to capture, since these are not directly related to the tradeoff between space, locality, and read efficiency.

Specifically, for databases of size  $N$ , our scheme offers both optimal space and optimal locality (i.e., space  $O(N)$  and locality  $O(1)$ ), and comes very close to offering optimal read efficiency as well. The read efficiency of our scheme when querying for a keyword  $w$  depends on the length of the list  $\text{DB}(w)$  that is associated with  $w$  (that is, the read efficiency depends on the number of identifiers that are associated with  $w$ ).<sup>3</sup> When querying for a keyword that is associated with  $n = N^{1-\epsilon(n)}$  identifiers, the read efficiency of our scheme is  $f(N) \cdot \epsilon(n)^{-1} + O(\log \log \log N)$ , where  $f(N) = \omega(1)$  may be any pre-determined function, and  $\omega(1) \cdot \epsilon(n)^{-1}$  is a lower bound as proved by Asharov et al. [ANS<sup>+</sup>16]. In particular, for any keyword that is associated with at most  $N^{1-1/o(\log \log \log N)}$  identifiers (i.e., for any keyword that is not exceptionally common), the read efficiency of our scheme when retrieving its identifiers is  $O(\log \log \log N)$ . We refer the reader to Sections 1.2 and 4 for a high-level overview and for a detailed description of this framework and of our new scheme, respectively.

Our results in the pad-and-split and statistical-independence frameworks, which are summarized in Table 1 and presented in more detail in Section 1.2, show a significant gap between the performance guarantees that can be offered within these two frameworks. In both frameworks we establish tight bounds that capture the basic techniques underlying all of the existing SSE schemes. Thus, any attempt to further improve upon the tradeoff between the space overhead, locality and read efficiency of our schemes must be based on new techniques that deviate from all known SSE schemes.

	Space	Locality	Read Efficiency
<b>This work (Thm. 1.1):</b> Pad-and-split lower bound	$\Omega(N \log N / \log L)$	$L$	$O(1)$
[DP17]: Pad-and-split scheme	$O(N \log N / \log L)$	$L$	$O(1)$
[ANS <sup>+</sup> 16]: Statistical-independence lower bound	$O(N)$	$O(1)$	$\omega(1) \cdot \epsilon(n)^{-1}$
[ANS <sup>+</sup> 16]: Statistical-independence scheme	$O(N)$	$O(1)$	$\tilde{O}(\log \log N)$
<b>This work (Thm. 1.2):</b> Statistical-independence scheme	$O(N)$	$O(1)$	$\omega(1) \cdot \epsilon(n)^{-1} + O(\log \log \log N)$

**Table 1: A summary of our contributions.** We denote by  $N$  the size of the database. The read efficiency in the lower bound of Asharov et al. [ANS<sup>+</sup>16] and in our statistical-independence scheme (Thm. 1.2) when querying for a keyword  $w$  depends on the number  $n = N^{1-\epsilon(n)}$  of identifiers that are associated with  $w$ .

In addition, our statistical-independence scheme is based on the modest assumption that no keyword is associated with more than  $N/\log^3 N$  identifiers, whereas the scheme of Asharov et al. [ANS<sup>+</sup>16] is based on the stronger assumption that no keyword is associated with more than  $N^{1-1/\log \log N}$  identifiers (thus, the read efficiency of their scheme does not contradict their lower bound, and our scheme has better read efficiency compared to their scheme). Finally, we note that the  $\omega(1)$  term in the read efficiency of our scheme can be set to any super-constant function (e.g.,  $\log \log \log N$ ).

## 1.2 Overview of Our Contributions

In this section we provide an overview of the two frameworks that we consider in this work, and present our results within each framework. As standard in the line of research on searchable sym-

<sup>3</sup>We emphasize that this does not hurt the security of SSE schemes, and still results in minimal leakage as required.

metric encryption, we represent a database 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 identifiers that are associated with each keyword  $w$ . We denote by  $N = \sum_{i=1}^{n_W} |\text{DB}(w_i)|$  the size of the database.

**Our pad-and-split framework.** Our pad-and-split framework considers schemes that are characterized by an algorithm denoted **SplitList** and consist of two phases. In the first phase, given a database  $\text{DB} = \{\text{DB}(w_1), \dots, \text{DB}(w_{n_W})\}$  of size  $N$ , for each keyword  $w_i$  the scheme invokes the **SplitList** algorithm on the length  $n_i$  of its corresponding list  $\text{DB}(w_i)$ , to obtain a vector  $(x_i^{(1)}, \dots, x_i^{(m)})$  of integers. The scheme then potentially pads the list  $\text{DB}(w_i)$  by adding “dummy” elements, and splits the padded list into sublists of lengths  $\text{len}^{(1)}, \dots, \text{len}^{(m)}$ , where  $x_i^{(j)}$  denotes the number of sublists of each length  $\text{len}^{(j)}$ . Then, in the second phase, for each possible length  $\text{len}^{(j)}$ , the scheme groups together all sublists of length  $\text{len}^{(j)}$ , and independently processes each such group to produce an encrypted database **EDB**.

We consider any possible instantiation of the **SplitList** algorithm (satisfying the necessary requirement that no list is longer than the sum of lengths of its sublists), and this enables us to describe a general template for constructing an SSE scheme based on any such algorithm given any one-way function. Our template yields schemes whose space usage and locality are essentially inherited from similar properties of their underlying **SplitList** algorithm, and whose read efficiency is always constant. We then demonstrate that this template captures the memory access patterns underlying essentially all existing schemes other than those of Asharov et al. [ANS<sup>+</sup>16]. Specifically, we show that each of these schemes can be obtained as an instantiation of our template using a suitable **SplitList** algorithm.

**A tight lower bound for pad-and-split schemes.** Equipped with our general notion of pad-and-split schemes, we prove a lower bound on the asymptotic efficiency guarantees of such schemes. Whereas the lower bound of Cash and Tessaro [CT14] states that SSE schemes with non-overlapping reads cannot simultaneously offer asymptotically-optimal space overhead and locality, we prove the following lower bound (capturing the same existing schemes) stating that the efficiency guarantees of pad-and-split schemes must in fact be very far from optimal:

**Theorem 1.1.** *Any pad-and-split SSE scheme for databases of size  $N$  with locality  $L = L(N)$  uses space  $\Omega(N \log N / \log L)$ .*

We show that this lower bound is tight, as it matches the tradeoff offered by the scheme of Demertzis and Papamanthou [DP17] (i.e., their scheme is an optimal instantiation of our framework). We refer the reader to Section 3 for a detailed and more formal presentation of our results, including an in-depth discussion of the existing pad-and-split instantiations.

**The statistical-independence framework.** The statistical-independence framework of Asharov et al. [ANS<sup>+</sup>16] considers symmetric searchable encryption schemes that are characterized by a pair of algorithms, denoted **RangesGen** and **Allocation**, and consist of two phases. In the first phase, given a database  $\text{DB} = \{\text{DB}(w_1), \dots, \text{DB}(w_{n_W})\}$  of size  $N$ , for each keyword  $w_i$  the scheme invokes the **RangesGen** algorithm on the length  $n_i$  of its corresponding list  $\text{DB}(w_i)$ , to obtain a set of *possible locations* in which the scheme may place the elements of the list  $\text{DB}(w_i)$ .<sup>4</sup> Then, in the second phase, given the sets of possible locations for all keywords, the scheme invokes the **Allocation** algorithm on

<sup>4</sup>Looking ahead, when supplied with a token corresponding to a keyword  $w_i$ , the server will return to the client all data stored in the possible locations of the list  $\text{DB}(w_i)$  (the server will not actually know in which of the possible locations the elements of the list are actually placed).

these sets to obtain the *actual locations* for the corresponding lists. A key property of this framework is that the **RangesGen** algorithm, which determines the set of possible locations for each list  $\text{DB}(w_i)$ , is applied separately and independently to the length of each list. Thus, the possible locations of each list are independent of the possible locations of all other lists (in contrast, the actual locations of the lists are naturally correlated).

Asharov et al. referred to a pair (**RangesGen**, **Allocation**) of such algorithms as an allocation scheme, and showed that any such allocation scheme can be used to construct an SSE scheme. Then, by constructing two allocation schemes they obtained two SSE schemes with space  $O(N)$  and locality  $O(1)$ . Without making any assumptions on the structure of the database, their first scheme has read efficiency  $\tilde{O}(\log N)$ , and under the assumption that no keyword is associated with more than  $N^{1-1/\log \log N}$  identifiers, their second scheme has read efficiency  $\tilde{O}(\log \log N)$ .

**Our leveled two-choice scheme.** Within the statistical-independence framework, as discussed above, we construct a scheme whose tradeoff between space, locality, and read efficiency matches the lower bound proved by Asharov et al. for scheme in this framework to within an additive  $O(\log \log \log N)$  factor in its read efficiency (see Section 4 for a formal statement of their lower bound).

Specifically, we construct a scheme whose read efficiency when querying for a keyword  $w$  depends on the length of the list  $\text{DB}(w)$  that is associated with  $w$  (that is, the read efficiency depends on the number of identifiers that are associated with  $w$ ). For any  $n \leq N$  we denote by  $r(N, n)$  the read efficiency when retrieving a list of length  $n$ , and prove the following theorem:

**Theorem 1.2.** *Assuming the existence of any one-way function, for any function  $f(N) = \omega(1)$  there exists an adaptively-secure symmetric searchable encryption scheme for databases of size  $N$  in which no keyword is associated with more than  $N/\log^3 N$  identifiers, with the following guarantees:*

- *Space  $O(N)$ .*
- *Locality  $O(1)$ .*
- *Read efficiency  $r(N, n) = f(N) \cdot \epsilon(n)^{-1} + O(\log \log \log N)$ , where  $n = N^{1-\epsilon(n)}$ .*
- *Token size  $O(1)$ .*

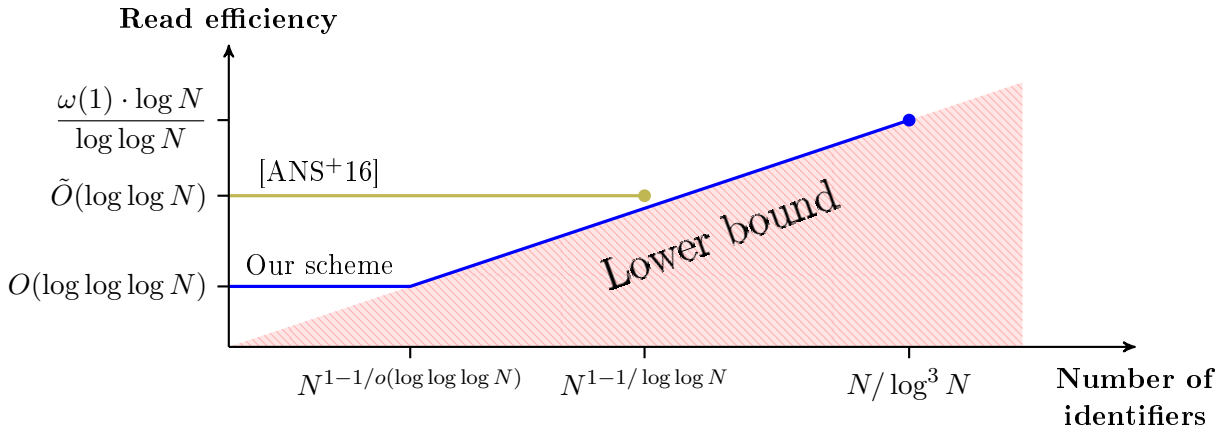
Our construction applies to databases of size  $N$  under the modest assumption that no keyword is associated with more than  $N/\log^3 N$  identifiers (note that the construction of Asharov et al. [ANS<sup>+</sup>16] is based on the stronger assumption that no keyword is associated with more than  $N^{1-1/\log \log N}$  identifiers). One can always generically deal (in a secure manner) with such extremely-common keywords by first excluding them from the database and applying our proposed scheme, and then applying in addition any other scheme for these extremely-common keywords (e.g., the “one-choice scheme” of Asharov et al. [ANS<sup>+</sup>16] or the recent scheme of Demertzis, Papadopoulos and Papamanthou [DPP17] – see Section 1.3 for more details).

When comparing our scheme to the scheme of Asharov et al. (see Table 1), both schemes offer space  $O(N)$  and locality  $O(1)$ , where the read efficiency of our scheme is strictly better than the read efficiency of their scheme – see Figure 1. In particular, for any keyword that is not exceptionally frequent (specifically, associated with at most  $N^{1-1/o(\log \log \log N)}$  identifiers), our scheme provides read efficiency  $O(\log \log \log N)$  whereas their scheme provides read efficiency  $\tilde{O}(\log \log N)$ .

**The structure of our scheme.** Our scheme is a *leveled* generalization of the “two-choice” scheme of Asharov et al. and consists of three levels for storing the elements of a given database. The first level consists of the two-choice SSE scheme of Asharov et al. but with *an exponentially improved read*

*efficiency.* Our key observation is that when viewing the first level as a collection of “bins”, then by allowing a few elements to “overflow” we can reduce the maximal load of each bin from  $\tilde{O}(\log \log N)$  (as in [ANS<sup>+</sup>16]) to  $O(\log \log \log N)$  and also handle much longer lists (i.e., much more frequent keywords). This then translates into improving the read efficiency in this level from  $\tilde{O}(\log \log N)$  to  $O(\log \log \log N)$ , while still using space  $O(N)$  and locality  $O(1)$ .

At this point, however, we have to store the overflowing elements. We store the vast majority of these elements in our second level, which consists of roughly  $\log N$  cuckoo hashing tables [PR04], where the  $j$  hash table is designed to store at most  $\hat{N}/2^j$  values each of which of size  $2^j$ . Our specific choice of cuckoo hashing as a static dictionary (i.e., a hash table) is due to its specific properties that guarantee the security of our scheme (see Section 2.4 for a discussion of these specific properties). In particular, our third level consists of a cuckoo hashing *stash* for each of the second-level cuckoo hashing tables. The goal of introducing this level is to reduce the failure probability of cuckoo hashing from noticeable to negligible, which is essential for the security of our resulting SSE scheme. We refer the reader to Section 4 for a detailed description of our scheme.



**Figure 1: The read efficiency of our statistical-independence scheme compared to that of Asharov et al. [ANS<sup>+</sup>16] and to the lower bound.** The read efficiency of our scheme is depicted by the blue line, and the read efficiency of the scheme of Asharov et al. is depicted by the yellow line (recall that our scheme supports keywords that are associated with up to  $N/\log^3 N$  identifiers, whereas the scheme of Asharov et al. only supports keywords that are associated with at most  $N^{1-1/\log \log N}$  identifiers). The read efficiency lower bound of Asharov et al. is depicted by the red triangle (note that it coincides with our blue line for keywords that are associated with at least  $N^{1-1/o(\log \log \log N)}$  and at most  $N/\log^3 N$  identifiers). In all three cases the read efficiency is presented as a function of the number of identifiers that are associated with the queried keyword.

**Directions for future research.** In this work we establish tight bounds on the tradeoff between the space overhead, locality and read efficiency of SSE schemes within two general frameworks. Although these two frameworks capture the memory access pattern underlying all existing schemes, there is clearly no guarantee that our lower bounds cannot be circumvented by following other approaches. Thus, the main open problem that arises from our work is to prove tight bounds for all SSE schemes. Additional natural open problems are to prove such tight bounds for dynamic SSE schemes, and to study the above tradeoff for searchable encryption in the public-key setting.

### 1.3 Related Work

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 functionality for SSE, as well as the first 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 searchable symmetric encryption schemes with various efficiency properties, support for data updates, authenticity, support for more advanced searches, and more (see [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, ANS<sup>+</sup>16, DP17] and the references therein). The two frameworks that we consider in this work capture schemes that satisfy that standard notions of SSE introduced by Curtmola et al. [CGK<sup>+</sup>06, CGK<sup>+</sup>11]. These schemes are discussed in Section 3.2 as instantiations of our pad-and-split framework, and in Section 4.2 as instantiations of the statistical-independence framework of Asharov et al. [ANS<sup>+</sup>16].

Our statistical-independence scheme can be applied to any database in which no keyword is associated with more than  $N/\log^3 N$  identifiers. As discussed above, one can always generically deal (in a secure manner) with such extremely-frequent keywords by first excluding them from the database and applying our proposed scheme, and then applying in addition any other scheme for these extremely-common keywords. For example, for these keywords one can apply the “one-choice scheme” of Asharov et al. or the recent scheme of Demertzis, Papadopoulos and Papamanthou [DPP17] that provides a sub-logarithmic read efficiency when searching for extremely frequent keywords<sup>5</sup>. Specifically, Demertzis et al. proposed a scheme that handles such extremely frequent keywords and improves their read efficiency from  $\tilde{O}(\log N)$  as guaranteed by the “one-choice scheme” of Asharov et al. to  $O(\log^{2/3+\delta} N)$  for any fixed constant  $\delta > 0$  (for all other keywords they use the two schemes of Asharov et al., which can now be replaced by our new scheme in its appropriate range of parameters).

## 1.4 Paper Organization

The remainder of this paper is organized as follows. In Section 2 we review the standard notion of symmetric searchable encryption schemes, as well as various tools that are used in our constructions. Then, in Section 3 we put forward our pad-and-split framework and then present our lower bound and new scheme in this framework. In Section 4 we review the statistical-independence framework and then present our new scheme in this framework.

## 2 Preliminaries

In this section we present the notation, definitions, and basic tools 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. For a probabilistic algorithm  $\text{Alg}$ , we denote its output when using fresh uniform random tape by  $\text{output} \leftarrow \text{Alg}(\text{input})$ . Additionally, we denote its (deterministic) output when using an explicit random tape  $r$  by  $\text{Alg}(\text{input}; r)$ .

---

<sup>5</sup>The scheme of Demertzis, Papadopoulos and Papamanthou [DPP17] is not captured by the two frameworks we consider in this work, as it requires the server to modify its stored data (i.e., the encrypted database) and the user to update her local state whenever a search query is issued.

## 2.1 Searchable Symmetric Encryption

Let  $W = \{w_1, \dots, w_{n_W}\}$  denote a set of keywords, where each keyword  $w_i$  is associated with a list  $DB(w_i) = \{\text{id}_1, \dots, \text{id}_{n_i}\}$  of document identifiers (these may correspond, for example, to documents in which the keyword  $w_i$  appears). A database  $DB = \{DB(w_1), \dots, DB(w_{n_W})\}$  consists of several such lists. We assume that each keyword and document identifier can be represented using a constant number of machine words, each of length  $O(\lambda)$  bits, in the unit-cost RAM model<sup>6</sup>. There are various different syntaxes for SSE schemes in the literature, where the main differences are in the flavor of interaction between the server and the client with each query. In this work we consider both a setting where the server decrypts the set of identifiers by itself, and a setting where the server does not decrypt this but rather sends encrypted data back to the client (who can then decrypt and learn the set of identifiers).

### 2.1.1 Functionality

A searchable symmetric encryption scheme is a 5-tuple  $(\text{KeyGen}, \text{EDBSetup}, \text{TokGen}, \text{Search}, \text{Resolve})$  of probabilistic polynomial-time algorithms satisfying the following requirements:

- The key-generation algorithm **KeyGen** takes as input the security parameter  $\lambda \in \mathbb{N}$  in unary representation and outputs a secret key  $K$ .
- The database setup **EDBSetup** algorithm takes as input a secret key  $K$  and a database  $DB$ , and outputs an encrypted database  $EDB$ .
- The token-generation algorithm **TokGen** takes as input a secret key  $K$  and a keyword  $w$ , and outputs a token  $\tau$  and some internal state  $\rho$ .
- The search algorithm **Search** takes as input a token  $\tau$  and an encrypted database  $EDB$ , and outputs a list  $R$  of results.
- The resolve algorithm **Resolve** takes as input a list  $R$  of results and an internal state  $\rho$ , and outputs a list  $M$  of document identifiers.

An SSE scheme for databases of size  $N = N(\lambda)$  is correct if for any database  $DB$  of size  $N$  and for any keyword  $w$ , with an overwhelming probability in the security parameter  $\lambda \in \mathbb{N}$ , it holds that  $M = DB(w)$  at the end of the following experiment:

1.  $K \leftarrow \text{KeyGen}(1^\lambda)$ .
2.  $EDB \leftarrow \text{EDBSetup}(K, DB)$ .
3.  $(\tau, \rho) \leftarrow \text{TokGen}(K, w)$ .
4.  $R \leftarrow \text{Search}(\tau, EDB)$ .
5.  $M = \text{Resolve}(\rho, R)$ .

We note that one can also consider a more adversarially-flavored notion of correctness, where an adversary adaptively interacts with a server with the goal of producing a query that results in an incorrect output. We refer the reader to [ANS<sup>+</sup>16] for more details, and here we only point out that our schemes in this paper satisfy such a notion as well.

### 2.1.2 Efficiency Measures

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

---

<sup>6</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).

**Space.** A symmetric searchable encryption scheme  $(\text{KeyGen}, \text{EDBSetup}, \text{TokGen}, \text{Search}, \text{Resolve})$  uses space  $s = s(\lambda, N)$  if for any  $\lambda, N \in \mathbb{N}$ , for any database  $\text{DB}$  of size  $N$ , and for any key  $K$  produced by  $\text{KeyGen}(1^\lambda)$ , the algorithm  $\text{EDBSetup}(K, \text{DB})$  produces encrypted databases that can be represented using  $s$  machine words.

**Locality.** The search procedure of any SSE scheme can be decomposed into a sequence of contiguous reads from the encrypted database  $\text{EDB}$ , and the locality is defined as the number of such reads. Specifically, locality is defined by viewing the  $\text{Search}$  algorithm of an SSE scheme as an algorithm that does not obtain as input the actual encrypted database, but rather only obtains oracle access to it. Each query to this oracle consists of an interval  $[a_i, b_i]$ , and the oracle replies with the machine words that are stored in this interval of  $\text{EDB}$ . At first, the  $\text{Search}$  algorithm is invoked on a token  $\tau$  and queries its oracle with some interval  $[a_1, b_1]$ . Then, it iteratively continues to compute the next interval to read based on  $\tau$  and all previously read intervals. 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.

**Read efficiency.** The notion of read efficiency compares the overall size of the portion of  $\text{EDB}$  that is read on each query to the size of the actual answer to the query. 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 Security Notions

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  $\boldsymbol{\tau} = \{\tau_i\}_{w_i \in \mathbf{w}}$ , 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

$\text{EDB}, \boldsymbol{\tau}$  (i.e.,  $(\text{EDB}, \boldsymbol{\tau}) \leftarrow \mathcal{S}(\mathcal{L}(\text{DB}, \mathbf{w}))$ ). The execution follows by giving  $(\text{EDB}, \boldsymbol{\tau})$  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 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 information 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

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

**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

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

## 2.2 Private-Key Encryption

We rely on the standard notion of a private-key encryption scheme with pseudorandom ciphertexts, which is easily realized based on the minimal assumption that a one-way function exists [Gol04]. For simplicity, and without loss of generality, we consider private-key encryption schemes in which the key-generation algorithm produces a  $\lambda$ -bit uniformly-distributed key, where  $\lambda \in \mathbb{N}$  is the security parameter. Thus, such an encryption scheme is fully specified via its encryption and decryption algorithms, which we denote by **Enc** and **Dec**, respectively.

**Definition 2.5.** Let  $\Pi = (\text{Enc}, \text{Dec})$  be a private-key encryption scheme. Then  $\Pi$  has *pseudorandom ciphertexts* if for every probabilistic polynomial-time algorithm  $\mathcal{A}$  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  $\text{Enc}_K(\cdot)$ ), and the above probabilities are taken over the choice of  $K \leftarrow \{0, 1\}^\lambda$ , the internal randomness of the algorithm **Enc** and the oracle  $\mathcal{R}$ .

## 2.3 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 (**HTSetup**, **HTLookup**). The algorithm **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 **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).

## 2.4 Cuckoo Hashing with a Stash

Cuckoo hashing is an efficient and practical hash table designed by Pagh and Rodler [PR04], providing worst-case constant lookup time and uses linear space. An important property of cuckoo hashing is that by storing a few elements in a secondary (small) data structure, referred to as a “stash”, it is possible to decrease its failure probability from noticeable to negligible [KMW09]. For our purposes in this work, it suffices to consider the following abstraction of cuckoo hashing with a stash:

- The memory is an abstract array  $[m]$ , where each cell may contain a single element or **NULL**.
- The potential locations of any element are randomly sampled (instead of being determined by hash functions).

We now summarize the abstract properties of cuckoo hashing with a stash in which we are interested for our construction in Section 4:

1. For storing  $n$  lists, where each list consists of  $\ell$  elements, an array of size  $O(n \cdot \ell)$  is used. The array is partitioned into two segments – a cuckoo hashing segment of size  $O(n \cdot \ell)$  and a stash segment of size  $s \cdot \ell$ .

2. Fetching a list requires accessing two random locations (of size  $\ell$  each) in the cuckoo hashing segment and accessing the entire stash segment.
3. When using a stash of size  $s = n^{o(1)}$ , the probability that  $n$  lists can be successfully stored is  $1 - O(n^{s/2})$  [ADW14, Thm. 2].<sup>7</sup>

### 3 The Pad-and-Split Framework: A Stronger Lower Bound

In this section we first formalize our pad-and-split framework for the design of symmetric searchable encryption schemes (Section 3.1). Then, we show that it captures the memory access patterns underlying essentially all of the existing symmetric searchable encryption schemes other than the schemes of Asharov et al. [ANS<sup>+</sup>16] (Section 3.2), and discuss the instantiation of Demertzis and Papamanthou (Section 3.3) whose tradeoff matches our lower bound (Section 3.4).

#### 3.1 The Pad-and-Split Framework

Our framework considers symmetric searchable encryption schemes that are characterized by a deterministic algorithm denoted **SplitList**, and consist of the following two phases:

- Given a database  $\text{DB} = \{\text{DB}(w_1), \dots, \text{DB}(w_{n_w})\}$  of size  $N$ , for each keyword  $w_i$  the scheme invokes the **SplitList** algorithm on the length  $n_i$  of its corresponding list  $\text{DB}(w_i)$ , to obtain a vector  $(x_i^{(1)}, \dots, x_i^{(m)})$  of integers. The scheme then potentially pads the list  $\text{DB}(w_i)$  by adding “dummy” elements, and splits the padded list into sublists of lengths  $\text{len}^{(1)}, \dots, \text{len}^{(m)}$ , where  $x_i^{(j)}$  denotes the number of sublists of each length  $\text{len}^{(j)}$ .
- For each possible length  $\text{len}^{(j)}$ , the scheme groups together all sublists of length  $\text{len}^{(j)}$ , and independently processes each such group to produce an encrypted database **EDB**.

A key property of our framework is that the **SplitList** algorithm, which determines the number of sublists of each length, does not take as input an actual list  $\text{DB}(w_i)$  but only its length  $n_i = |\text{DB}(w_i)|$ . This algorithm is parameterized by the possible lengths  $\text{len}^{(1)}, \dots, \text{len}^{(m)}$  of sublists, and also by upper bounds  $s^{(1)}, \dots, s^{(m)}$  on the total number of sublists of lengths  $\text{len}^{(1)}, \dots, \text{len}^{(m)}$ , respectively. We allow the parameters  $m, \text{len}^{(1)}, \dots, \text{len}^{(m)}$  and  $s^{(1)}, \dots, s^{(m)}$  to depend on the total length  $N = \sum_{i=1}^{n_w} |\text{DB}(w_i)|$  of the database, but do not explicitly denote this for ease of notation.

We consider any possible instantiation of the **SplitList** algorithm subject to satisfying two natural requirements. First, we require that each list  $\text{DB}(w_i)$  is split into sublists whose total length is at least the length of  $\text{DB}(w_i)$ . Second, we require that for every possible sublist length  $\text{len}^{(j)}$  there are at most  $s^{(j)}$  sublists of length  $\text{len}^{(j)}$  in the worst-case over all possible databases of size  $N$ . Formally:

**Definition 3.1.** We say that a **SplitList** algorithm, parameterized by  $(\text{len}^{(1)}, \dots, \text{len}^{(m)})$  and  $(s_1, \dots, s_m)$  is *valid* if for every integer  $N$  and for every vector of lengths  $(n_1, \dots, n_k)$  with  $\sum_{i=1}^k n_i = N$ , it holds that:

- **Each list is not longer than the sum of lengths of its sublists:** For every  $n_i$  it holds that  $x_i^{(1)} \cdot \text{len}^{(1)} + \dots + x_i^{(m)} \cdot \text{len}^{(m)} \geq n_i$ , where  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$ .
- **Each  $s^{(j)}$  upper bounds the number of sublists of length  $\text{len}^{(j)}$ :** For every  $j \in [m]$  it holds that  $\sum_{i=1}^k x_i^{(j)} \leq s^{(j)}$ , where  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$  for every  $i \in [k]$ .

<sup>7</sup>Note that in the original work of Kirsch et al. [KMW09] they considered a constant-sized stash, whereas in this work we are interested in a stash whose size is not necessarily constant, and thus we rely on [ADW14].

In addition, we say that `SplitList` has locality  $L$  if each list  $\text{DB}(w_i)$  is split into at most  $L$  sublists. Formally:

**Definition 3.2.** We say that a `SplitList` algorithm *has locality*  $L = L(N)$  if for every  $n_i$  it holds that  $x_i^{(1)} + \dots + x_i^{(m)} \leq L$ , where  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$ .

Equipped with our notion of a valid `SplitList` algorithm, we describe a general template (see Construction 3.5) for constructing symmetric searchable encryption schemes given any such algorithm. We rely, in addition, on a pseudorandom function PRF and a private-key encryption scheme  $(\text{Enc}, \text{Dec})$  with pseudorandom ciphertexts – both of which can be constructed based on any one-way function as discussed in Section 2.2. This yields the following theorem:

**Theorem 3.3.** *Given a valid `SplitList` algorithm with parameters  $(\text{len}^{(1)}, \dots, \text{len}^{(m)})$  and  $(s^{(1)}, \dots, s^{(m)})$ , a pseudorandom function PRF, and a private-key encryption scheme  $(\text{Enc}, \text{Dec})$  with pseudorandom ciphertexts, Construction 3.5 is a static  $\mathcal{L}_{\min}$ -secure symmetric searchable encryption scheme for databases of size  $N$  with the following parameters:*

- Space  $O\left(\sum_{j=1}^m s^{(j)} \cdot \text{len}^{(j)}\right)$ .
- Locality  $O(L(N))$ , where `SplitList` has locality  $L(N)$ .
- Read efficiency  $O(1)$ .
- Token size  $O(1)$ .

Moreover, Construction 3.5 is an adaptive  $\mathcal{L}_{\min}^{\text{adap}}$ -secure symmetric searchable encryption scheme in the random-oracle model, when instantiating PRF and  $(\text{Enc}, \text{Dec})$  appropriately.

Note that Theorem 3.3 guarantees that Construction 3.5 is statically-secure in the standard model (although, we do prove it can be made adaptively secure in the random-oracle model). The next theorem shows that a simple modification of the scheme (described as Construction 3.6), based on an idea sketched by Stefanov et al. [SPS14], is in fact adaptively secure in the standard model. This comes at the cost of increasing the token size from tokens of size  $O(1)$  to tokens of size  $O(L)$ , where  $L$  is the locality of the `SplitList` algorithm.<sup>8</sup> In this scheme, the client decrypts the results sent by the server (using the `Resolve` algorithm), and thus the scheme leaks only the size of the results. This is in contrast to the scheme described in Construction 3.5, where the server decrypts the results, and thus the scheme leaks the results themselves<sup>9</sup>.

**Theorem 3.4.** *Given a valid `SplitList` algorithm with parameters  $(\text{len}^{(1)}, \dots, \text{len}^{(m)})$  and  $(s^{(1)}, \dots, s^{(m)})$ , a pseudorandom function PRF, and a private-key encryption scheme  $(\text{Enc}, \text{Dec})$  with pseudorandom ciphertexts, Construction 3.6 is an adaptive  $\mathcal{L}_{\text{size}}^{\text{adap}}$ -secure symmetric searchable encryption scheme for databases of size  $N$  with the following parameters:*

- Space  $O\left(\sum_{j=1}^m s^{(j)} \cdot \text{len}^{(j)}\right)$ .
- Locality  $O(L(N))$ , where `SplitList` has locality  $L(N)$ .
- Read efficiency  $O(1)$ .
- Token size  $O(L(N))$ .

In the remainder of this section we first provide a high-level overview of these schemes, and then formally prove Theorems 3.3 and 3.4.

<sup>8</sup>We refer the reader to the work of Chase and Kamara [CK10] for a discussion on the necessity of using a random oracle for adaptive security with succinct search tokens.

<sup>9</sup>Note that any scheme in which the server decrypts the results can be easily transformed into a scheme where only the client decrypts the results by adding an additional encryption layer – but this does not necessarily hold in the other direction.

**Overview of the schemes.** In both schemes each list of document identifiers  $\text{DB}(w_i)$  is padded and split as dictated by the output  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, |\text{DB}(w_i)|)$ . That is,  $\text{DB}(w_i)$  is padded to length  $x_i^{(1)} \cdot \text{len}^{(1)} + \dots + x_i^{(m)} \cdot \text{len}^{(m)}$ , and split into sublists, where for each  $j \in [m]$  there are  $x_i^{(j)}$  sublists of length  $\text{len}^{(j)}$ . Then, we construct an encrypted database which consists of the following hash tables:

- A hash table that stores (in encrypted manner) the lengths of all lists, and is padded to contain exactly  $N$  elements.
- For every  $j \in [m]$  a hash table that stores (in an encrypted manner) all sublists of length  $\text{len}^{(j)}$ , and is padded to contain exactly  $s^{(j)}$  sublists of this length.

In all hash tables we store the various elements according to pseudorandom labels that are derived from each corresponding keyword  $w$  via a pseudorandom function whose key is known only to the client. Intuitively speaking, the scheme is secure for any valid **SplitList** algorithm due to the following three reasons: (1) The number of padded elements and the number of sublists each list is split into depend only on the length of each list, (2) each hash table consists of encrypted elements with pseudorandom labels, and (3) the size of each hash table depends only on the size of the database.

The main differences between Construction 3.5 (providing static security) and Construction 3.6 (providing adaptive security) are as follows:

1. The lengths of the lists in Construction 3.6 are encrypted using one-time pads. This is required in order to allow “explaining” a random value as the encryption of any particular length on the fly (given that adversaries may be adaptive).
2. In Construction 3.5, for each searched keyword the server is given keys derived from that keyword, allowing it to compute the labels and decrypt the document identifiers associated with that keyword. In Construction 3.6 the server is given the labels themselves (thus, the token size is  $O(L)$ ), and can only locate the encrypted document identifiers, but not to decrypt them.

**Proof of Theorem 3.3.** It is easy to see that the scheme is correct unless the same label appears more than once, which happens with at most a negligible probability. Also, it is easy to verify the space overhead, locality, read efficiency and token size of the construction.

For proving the security of the scheme, recall that in the real execution the adversary  $\mathcal{A}$  outputs  $\text{DB}$  and  $\mathbf{w}$ , 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  $\tau = \{\tau_i\}_{w_i \in \mathbf{w}}$  in such a way that the adversary cannot distinguish between these two executions. Consider the following simulator  $\mathcal{S}$ :

- **Input:**  $\mathcal{L}(\text{DB}, \mathbf{w}) = (N, \{\text{DB}(w_i)\}_{w_i \in \mathbf{w}})$ .
- **The simulator:**
  1. The simulator  $\mathcal{S}$  samples uniform random keys  $\tau_i = (\text{label}_i, K_i, \widehat{K}_i)$  for every  $w_i \in \mathbf{w}$ .
  2.  $\mathcal{S}$  initializes the sets  $T_1, \dots, T_m$  and  $T$  to be empty sets. For every  $w_i \in \mathbf{w}$  and given  $\text{DB}(w_i)$ , it computes  $n_i$ , invokes  $(x_1, \dots, x_m) \leftarrow \text{SplitList}(N, n_i)$ , follows Step 2d in the construction, and adds the pairs  $(\text{label}_{j,x}, d_{j,x})$  to the corresponding sets  $T_j$ . It also encrypts  $\widehat{n}_i = \text{Enc}_{K_i}(n_i)$  and adds the pair  $(\text{label}_i, \widehat{n}_i)$  to the set  $T$ .
  3.  $\mathcal{S}$  pads the set  $T$  to contain exactly  $N$  elements by adding dummy elements, and pads each set  $T_j$  to contain exactly  $s^{(j)}$  elements by adding dummy elements.



**CONSTRUCTION 3.5** (One-Round Pad-and-Split Symmetric Searchable Encryption).

A pad-and-split SSE scheme is parameterized by a `SplitList` algorithm, and by the following values (all values are functions of the size  $N$  of the database):

1. Locality parameter  $L$ .
2. Possible lengths  $\text{len}^{(1)}, \dots, \text{len}^{(m)}$  of sublists.
3. Upper bounds  $s^{(1)}, \dots, s^{(m)}$  on the total number of sublists of lengths  $\text{len}^{(1)}, \dots, \text{len}^{(m)}$ , respectively.

**Key generator.** The algorithm `KeyGen` on input  $1^\lambda$  samples and outputs a key  $K \leftarrow \{0, 1\}^\lambda$  for PRF.

**Setup.** The algorithm `EDBSetup` on input  $(K, \text{DB})$  is defined as follows:

1. Initialize  $t + 1$  empty sets  $T, T_1, \dots, T_m$ , where  $T$  will consist of the lengths of the lists, and each set  $T_j$  will consist of all sublists of length  $\text{len}^{(j)}$ .
2. For every keyword  $w_i \in W$  with an associated list  $\text{DB}(w_i) = \{\text{id}_1, \dots, \text{id}_{n_i}\}$ :
  - (a) Compute  $(\text{label}_i, K_i, \widehat{K}_i) = \text{PRF}_K(w_i)$ .
  - (b) Compute  $\widehat{n}_i = \text{Enc}_{K_i}(n_i)$  and add the pair  $(\text{label}_i, \widehat{n}_i)$  to the set  $T$ .
  - (c) Compute  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$ .
  - (d) For every  $j = 1, \dots, m$ :
    - i. For every  $x = 1, \dots, x_i^{(j)}$ :
      - A. Take the next  $\text{len}^{(j)}$  elements from the list  $\text{DB}(w_i)$  and create a block  $\{\text{id}'_1, \dots, \text{id}'_{\text{len}^{(j)}}\}$ . If there are less than  $\text{len}^{(j)}$  elements left in  $\text{DB}(w_i)$ , then pad with dummy elements.
      - B. Compute a label:  $\text{label}_{j,x} = \text{PRF}_{\widehat{K}_i}(j, x)$ .
      - C. Encrypt  $d_{j,x} = (\text{Enc}_{K_i}(\text{id}'_1), \dots, \text{Enc}_{K_i}(\text{id}'_{\text{len}^{(j)}}))$ .
      - D. Insert the pair  $(\text{label}_{j,x}, d_{j,x})$  into the set  $T_j$ .
3. Pad the set  $T$  to contain exactly  $N$  elements by adding dummy elements, and pad each set  $T_j$  to contain exactly  $s^{(j)}$  elements by adding dummy elements.
4. For each set  $T, T_1, \dots, T_m$ , uniformly shuffle the set, and generate a hash table by invoking the `HTSetup` algorithm for obtaining hash tables  $\text{HT}(T), \text{HT}(T_1), \dots, \text{HT}(T_m)$ .
5. Output  $\text{EDB} = (\text{HT}(T), (\text{HT}(T_1), \dots, \text{HT}(T_m)))$ .

**Token generator.** The algorithm `TokGen` on input  $(K, w_i)$  computes and outputs the token  $\tau_i = (\text{label}_i, K_i, \widehat{K}_i) = \text{PRF}_K(w_i)$ .

**Search.** The algorithm `Search` on input  $(\tau_i, \text{EDB})$ , where  $\tau_i = (\text{label}_i, K_i, \widehat{K}_i)$  and  $\text{EDB} = (\text{HT}(T), \text{HT}(T_1), \dots, \text{HT}(T_m))$ , is defined as follows:

1. Initialize a list of document identifiers  $R = \emptyset$ .
2. Invoke `HTLookup` on the hash table  $\text{HT}(T)$  and label  $\text{label}_i$  to retrieve  $\widehat{n}_i = \text{Enc}_{K_i}(n_i)$ . Decrypt  $\widehat{n}_i$  using the key  $K_i$ , and compute  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$ .
3. For every  $j \in [m]$  and for every  $x \in [x_i^{(j)}]$  compute  $\text{label}_{j,x} = \text{PRF}_{\widehat{K}_i}(j, x)$ . Invoke `HTLookup` on the hash table  $\text{HT}(T_j)$  for the label  $\text{label}_{j,x}$ , and obtain the block  $d_{j,x}$ . Decrypt the block using the key  $K_i$  and add the elements to the list  $R$ . For a block that contains dummy elements, obtain and decrypt only the part that does not contain dummy elements.

**CONSTRUCTION 3.6** (Two-Round Pad-and-Split Symmetric Searchable Encryption).

A pad-and-split SSE scheme is parameterized by a **SplitList** algorithm, and the following values (all values are functions of the size  $N$  of the database):

1. Locality parameter  $L$ .
2. Possible lengths  $\text{len}^{(1)}, \dots, \text{len}^{(m)}$  of sublists.
3. Upper bounds  $s^{(1)}, \dots, s^{(m)}$  on the total number of sublists of lengths  $\text{len}^{(1)}, \dots, \text{len}^{(m)}$ , respectively.

**Key generator.** The algorithm **KeyGen** on input  $1^\lambda$  samples a key  $K \leftarrow \{0, 1\}^\lambda$  for PRF, samples a key  $\widehat{K} \leftarrow \{0, 1\}^\lambda$  for (Enc, Dec), and outputs  $(K, \widehat{K})$ .

**Setup.** The algorithm **EDBSetup** on input  $((K, \widehat{K}), \text{DB})$  is defined as follows:

1. Initialize  $t + 1$  empty sets  $T, T_1, \dots, T_m$ , where  $T$  will consist of the lengths of the lists, and each set  $T_j$  will consist of all sublists of length  $\text{len}^{(j)}$ .
2. For every keyword  $w_i \in \mathcal{W}$  with an associated list  $\text{DB}(w_i) = \{\text{id}_1, \dots, \text{id}_{n_i}\}$ :
  - (a) Compute  $((\text{label}_i, K_i), (\text{label}_{i,1}, \dots, \text{label}_{i,L})) = \text{PRF}_K(w_i)$ .
  - (b) Compute  $\widehat{n}_i = K_i \oplus n_i$  and add the pair  $(\text{label}_i, \widehat{n}_i)$  to the set  $T$ .
  - (c) Compute  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$ .
  - (d) For every  $j = 1, \dots, m$ :
    - i. For every  $x = 1, \dots, x_i^{(j)}$ :
      - A. Take the next  $\text{len}^{(j)}$  elements from the list  $\text{DB}(w_i)$  and create a block  $\{\text{id}'_1, \dots, \text{id}'_{\text{len}^{(j)}}\}$ . If there are less than  $\text{len}^{(j)}$  elements left in  $\text{DB}(w_i)$ , then pad with dummy elements.
      - B. Let  $\text{label}$  be the first unused label from  $(\text{label}_{i,1}, \dots, \text{label}_{i,L})$ .
      - C. Encrypt  $d_{j,x} = (\text{Enc}_{\widehat{K}}(\text{id}'_1), \dots, \text{Enc}_{\widehat{K}}(\text{id}'_{\text{len}^{(j)}}))$ .
      - D. Insert the pair  $(\text{label}, d_{j,x})$  into the set  $T_j$ .
3. Pad the set  $T$  to contain exactly  $N$  elements by adding dummy elements, and pad each set  $T_j$  to contain exactly  $s^{(j)}$  elements by adding dummy elements.
4. For each set  $T, T_1, \dots, T_m$ , uniformly shuffle the set, and generate a hash table by invoking the **HTSetup** algorithm for obtaining hash tables  $\text{HT}(T), \text{HT}(T_1), \dots, \text{HT}(T_m)$ .
5. Output  $\text{EDB} = (\text{HT}(T), (\text{HT}(T_1), \dots, \text{HT}(T_m)))$ .

**Token generator.** The algorithm **TokGen** on input  $((K, \widehat{K}), w_i)$  computes and outputs the token  $\tau_i = ((\text{label}_i, K_i), (\text{label}_{i,1}, \dots, \text{label}_{i,L})) = \text{PRF}_K(w_i)$ .

**Search.** The algorithm **Search** on input  $(\tau_i, \text{EDB})$ , where  $\tau_i = ((\text{label}_i, K_i), (\text{label}_{i,1}, \dots, \text{label}_{i,L}))$  and  $\text{EDB} = (\text{HT}(T), \text{HT}(T_1), \dots, \text{HT}(T_m))$ , is defined as follows:

1. Initialize a list of results  $R = \emptyset$ .
2. Invoke **HTLookup** on the hash table  $\text{HT}(T)$  and label  $\text{label}_i$  to retrieve  $\widehat{n}_i = K_i \oplus n_i$ . Decrypt  $n_i = K_i \oplus \widehat{n}_i$  and compute  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$ .
3. For every  $j \in [m]$  and for every  $x \in [x_i^{(j)}]$ , let  $\text{label}$  be the first unused label from  $(\text{label}_{i,1}, \dots, \text{label}_{i,L})$ . Invoke **HTLookup** on the hash table  $\text{HT}(T_j)$  for the label  $\text{label}_{j,x}$ , obtain the block  $d_{j,x}$ , and add its elements to the list  $R$ . For a block that contains dummy elements, obtain only the part that does not contain dummy elements.

**Resolve.** The algorithm **Resolve** on input  $((K, \widehat{K}), R)$  computes and outputs the identifiers  $M = \{\text{Dec}_{\widehat{K}}(c) : c \in R\}$ .

4.  $\mathcal{S}$  uniformly shuffles each of the sets  $T, T_1, \dots, T_m$ . It then generates a hash table from each of these sets using the algorithm  $\text{HTSetup}$ , and defines  $\text{EDB} = (\text{HT}(T), (\text{HT}(T_1), \dots, \text{HT}(T_m)))$ .
5.  $\mathcal{S}$  outputs  $\text{EDB}$  and  $\tau = \{\tau_i\}_{w_i \in \mathbf{w}}$ .

We now claim that the adversary cannot distinguish between the pair  $(\text{EDB}, \tau)$  that it receives in the real execution and the same pair in the ideal execution with a non-negligible advantage. Consider the following sequence of 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 **Hyb<sub>0</sub>** 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  $\tau_i = (\text{label}_i, K_i, \widehat{K}_i)$  are uniformly distributed.
- **Hyb<sub>2</sub>**. This experiment is obtained from **Hyb<sub>1</sub>** by considering the set of keywords consisting of all  $w_i \in W \setminus \mathbf{w}$ , and replacing each one of the corresponding functions  $\text{PRF}_{\widehat{K}_i}(\cdot)$  with independent truly random functions. Observe that, as a result, that the elements  $\text{label}_{j,x} = \text{PRF}_{\widehat{K}_i}(j, x)$  derived in Step 2(d)iB are distributed uniformly.
- **Hyb<sub>3</sub>**. This experiment is obtained from **Hyb<sub>2</sub>** as follows: For every  $w_i \in W \setminus \mathbf{w}$ , we replace the values  $\widehat{n}_i$  (in Step 2b of  $\text{EDBSetup}$ ) with independent and uniformly-distributed values of the appropriate length. In addition, for every  $w_i \in W \setminus \mathbf{w}$ ,  $j \in [m]$  and  $x \in [x_i^{(j)}]$  (where  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$ ) we replace the encrypted values  $d_{j,x}$  (from Step 2(d)iC) with uniformly-distributed values of the appropriate length.
- **Hyb<sub>4</sub>**. This is the ideal execution, where we run the simulator  $\mathcal{S}$  defined above.

We observe that **Hyb<sub>0</sub>** and **Hyb<sub>1</sub>**, as well as **Hyb<sub>1</sub>** and **Hyb<sub>2</sub>**, are computationally indistinguishable based on the security of the pseudorandom function  $\text{PRF}$ . In addition, **Hyb<sub>2</sub>** and **Hyb<sub>3</sub>** are computationally indistinguishable based on the pseudorandom ciphertexts property of the encryption scheme  $\text{Enc}$ . Finally, the experiments **Hyb<sub>3</sub>** and **Hyb<sub>4</sub>** are identical by the definition of the simulator  $\mathcal{S}$ .

To prove adaptive security in the random oracle model, we use similar technique as in [CT14] and [ANS<sup>+</sup>16]. We instantiate  $\text{PRF}_K(x) = H(K||x)$  and  $\text{Enc}_K(x; r) = (r, H(K||r) \oplus x)$ , where  $H$  is the random oracle. Then, for any derived  $\text{PRF}$  key  $\widehat{K}_i$ , label  $\text{label}$ , and indices  $j$  and  $x$ , we can program the random oracle such that  $\text{Enc}_{\widehat{K}_i}(x, j) = \text{label}$ , by fixing  $H(\widehat{K}_i||j, x) = c \oplus \text{label}$ . Also, for any derived encryption key  $K_i$ , random  $(r, c)$ , and a message  $m$ , we can program the random oracle such that  $\text{Enc}_{K_i}(x; r) = (r, c)$ , by fixing  $H(K_i||r) = c \oplus m$ .

For the security proof, the simulator in the first step outputs  $\text{EDB} = (\text{HT}(T), (\text{HT}(T_1), \dots, \text{HT}(T_m)))$ , where  $T$  and each  $T_i$  contains an appropriate amount of random pairs of the appropriate length. Later, upon each query  $w_i$  from the adversary, the simulator learns  $\text{DB}(w_i) = \{\text{id}_1, \dots, \text{id}_{n_i}\}$ . It then samples keys  $K_i$  and  $\widehat{K}_i$ , computes  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$ . It choose a random unused pair  $(\text{label}_i, (c, r))$  from  $T$ , and programs the oracle such that  $\text{Dec}_{K_i}(r, c) = n_i$ . Similarly, for each  $j \in [m]$ , the simulator chooses  $x_i^{(j)}$  random unused pairs from  $H_j$ , denoted by  $(\text{label}_{j,1}, d_{j,1}), \dots, (\text{label}_{j,x_i^{(j)}}, d_{j,x_i^{(j)}})$ , and programs the oracle such that for each  $x \in [x_i^{(j)}]$  it holds that  $\text{PRF}_{\widehat{K}_i}(j, x) = \text{label}_{j,x}$ , and that decrypting  $d_{j,x_i^{(j)}}$  using  $K_i$  results in the next  $\text{len}^{(j)}$  identifiers in  $\text{DB}(w_i)$ . If at any point in time the programming cannot succeed (because the corresponding values are already set for  $H$ ) then the simulator aborts. Otherwise, it outputs  $\tau_i = (\text{label}_i, K_i, \widehat{K}_i)$ . We omit the formal analysis as it follows standard techniques. ■

**Proof of Theorem 3.4.** It is easy to see that the scheme is correct unless the same label appears more than once, which happens with at most a negligible probability. Also, it is easy to verify the space overhead, locality, read efficiency and token size of the construction.

For proving the adaptive security of the scheme, recall that in the real execution the adversary  $\mathcal{A}$  outputs  $\text{DB}$ , and is given the encrypted database  $\text{EDB}$ . Then, it adaptively issues keyword queries, where a query  $w_i$  is answered by  $\tau_i = \text{TokGen}(K, w_i)$ . In the ideal execution, the simulator  $\mathcal{S}$  is given only  $N$  and should output  $\text{EDB}$ , and upon each query  $w_i$  the simulator is given  $|\text{DB}(w_i)|$  and should output a token  $\tau_i$ , in such a way that the adversary cannot distinguish between these two executions. Consider the following simulator  $\mathcal{S}$ :

- **Initialization phase.** The simulator receives the leakage  $\mathcal{L}_{\text{size}}^{\text{adap}}(\text{DB}) = N$ .
  1.  $\mathcal{S}$  initializes the set  $T$  to contain  $N$  random pairs of the appropriate length.
  2.  $\mathcal{S}$  initializes the sets  $T_1, \dots, T_m$ , where each set  $T_j$  contains exactly  $s^{(j)}$  random pairs of the appropriate length.
  3.  $\mathcal{S}$  generates a hash table from each set using the algorithm  $\text{HTSetup}$ , and outputs  $\text{EDB} = (\text{HT}(T), (\text{HT}(T_1), \dots, \text{HT}(T_m)))$ .
  4.  $\mathcal{S}$  stores the tables  $T_1, \dots, T_m$  and  $T$  in an internal state.
- **Query.** With each query that occurs, the simulator receives the leakage  $n_i = |\text{DB}(w_i)|$ .
  1.  $\mathcal{S}$  computes  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$ .
  2.  $\mathcal{S}$  chooses a random pair  $(\text{label}_i, r_i)$  from  $T$  and removes the pair.
  3. It computes  $K_i = n_i \oplus r_i$ .
  4. For each  $j \in [m]$ ,  $\mathcal{S}$  chooses  $x_i^{(j)}$  random pairs from  $T_j$ , remembers their labels, and removes the pairs.
  5. Let  $\text{label}_{i,1}, \dots, \text{label}_{i,\ell}$  be the labels of the removed pairs from the last step. It holds that  $\ell \leq L$ , so  $\mathcal{S}$  samples  $L - \ell$  additional labels  $\text{label}_{i,\ell+1}, \dots, \text{label}_{i,L}$ .
  6.  $\mathcal{S}$  outputs  $((\text{label}_i, K_i), (\text{label}_{i,1}, \dots, \text{label}_{i,L}))$ .

We now claim that the adversary cannot distinguish between the the real execution and ideal execution with a non-negligible advantage. Consider the following sequence of hybrid experiments:

- **Hyb<sub>0</sub>.** This is the real execution, where the adversary  $\mathcal{A}$  initially receives  $\text{EDB}$ , and upon each query  $w_i$  it receives  $\tau_i = \text{PRF}_K(w_i) = ((\text{label}_i, K_i), (\text{label}_{i,1}, \dots, \text{label}_{i,L}))$  for every queried  $w_i$ .
- **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  $\tau_i = ((\text{label}_i, K_i), (\text{label}_{i,1}, \dots, \text{label}_{i,L}))$  are uniformly distributed.
- **Hyb<sub>2</sub>.** This experiment is obtained from  $\text{Hyb}_1$  as follows: For every  $w_i \in W$ , we replace the values  $\hat{n}_i$  (in Step 2b of  $\text{EDBSetup}$ ) with independent and uniformly-distributed values  $r_i$  of the appropriate length. In addition, upon each query  $w_i$  we replace  $K_i$  in  $\tau_i$  with  $r_i \oplus n_i$ .
- **Hyb<sub>3</sub>.** This experiment is obtained from  $\text{Hyb}_2$  as follows: For every  $w_i \in W$ ,  $j \in [m]$  and  $x \in [x_i^{(j)}]$  (where  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$ ), we replace the encrypted values  $d_{i,x}$  (from Step 2(d)iC) with uniformly-distributed values of the appropriate length. Observe that, as a result, the key  $\hat{K}$  that is produced by  $\text{KeyGen}$  is redundant.
- **Hyb<sub>4</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  $\text{PRF}$ . In addition,  $\text{Hyb}_1$  and  $\text{Hyb}_2$  are identically distributed. Next,  $\text{Hyb}_2$

and  $\text{Hyb}_3$  are computationally indistinguishable based on the pseudorandom ciphertexts property of the encryption scheme  $\text{Enc}$ . Finally, the experiments  $\text{Hyb}_3$  and  $\text{Hyb}_4$  are identical by the definition of the simulator  $\mathcal{S}$ . ■

### 3.2 The Generality of the Pad-and-Split Framework

We now demonstrate that our pad-and-split framework captures the memory access patterns underlying the vast majority of existing symmetric searchable encryption schemes for supporting keywords search (i.e., we show that these schemes can be obtained as instantiations of our framework). We note that each of these schemes offers other important implementation details, improvements and optimizations that we do not intend to capture using our framework (since these are not directly related to the tradeoff between space, locality, and read efficiency), and we refer to the relevant papers for further details.

**The scheme of Curtmola et al. [CGK<sup>+</sup>06].** This is the most common technique underlying the vast majority of existing schemes (in particular, [CK10, vLSD<sup>+</sup>10, CGK<sup>+</sup>11, KPR12, CJJ<sup>+</sup>13, KO13]). In this scheme each list is split into single elements (i.e., sublists of length 1), and those are stored in the same hash table. This is captured by our framework when setting  $m = 1$ ,  $\text{len}^{(1)} = 1$ ,  $s^{(1)} = N$ , and  $\text{SplitList}(N, n_i) = (n_i)$ . This results in a scheme with space  $O(N)$ , locality  $O(N)$ , and read efficiency  $O(1)$ .

**The 2lev scheme of Cash et al. [CJJ<sup>+</sup>14].** This scheme can be viewed as a pad-and-split scheme with two possible lengths,  $b$  and  $B$ , where  $b < B$ . A list of length at most  $b$  is padded to length  $b$ , and a list of length greater than  $b$  is padded to a length that is a multiple of  $B$  and then split into sublists of length  $B$  (in order to reduce space overhead, this scheme does not add dummy lists, thus resulting in a non-standard leakage function). This results in a scheme with space  $O(N \cdot (b + \frac{B}{b+1}))$ , locality  $O(N/B)$ , and read efficiency  $O(1)$ .

**A simple scheme with  $O(N^2)$  space.** In this scheme each list is padded to the maximal possible length (i.e., to length  $N$ , where  $N = \sum_{i=1}^{n_w} |\text{DB}(w_i)|$ ), and all lists are stored in the same hash table. This is captured by our framework when setting  $m = 1$ ,  $\text{len}^{(1)} = N$ ,  $s^{(1)} = N$  and  $\text{SplitList}(N, n_i) = (1)$ . This results in a scheme with space  $O(N^2)$ , locality  $O(1)$ , and read efficiency  $O(1)$ .

**The scheme of Cash and Tessaro [CT14].** This scheme splits a list of length  $n_i$  into at most  $\log n_i$  sublists of lengths that are powers of 2 according to the binary representation of  $n_i$ . Then, for each possible power of 2, the scheme stores sublists of that length in a separate hash table. This is captured by our framework when setting  $m = \lceil \log N \rceil + 1$ ,  $\text{len}^{(j)} = 2^{j-1}$ ,  $s^{(j)} = N/2^{j-1}$ , and the  $\text{SplitList}$  algorithm on input  $n_i$  outputs a binary vector of length  $m$  which corresponds to the binary representation of  $n_i$ . This results in a scheme with space  $O(\sum_{j=1}^m \text{len}^{(j)} s^{(j)}) = O(N \log N)$ , locality  $O(\log N)$ , and read efficiency  $O(1)$ .

**The scheme of Asharov et al. [ANS<sup>+</sup>16, Sec. 5].** This scheme improves the one of Cash and Tessaro [CT14]. In this scheme, a list of length  $2^{p_i-1} < n_i \leq 2^{p_i}$  is padded to length  $2^{p_i}$  and stored as a whole. This is captured by our framework when setting  $m = \lceil \log N \rceil + 1$ ,  $\text{len}^{(j)} = 2^{j-1}$ ,  $s^{(j)} = 2N/2^{j-1}$ , and the  $\text{SplitList}$  algorithm, on input  $n_i$ , outputs a vector of length  $m$  where all the

entries are zeros except for a one that appears in the location  $\lceil \log n_i \rceil + 1$ . This results in a scheme with space  $O(\sum_{j=1}^m \text{len}^{(j)} s^{(j)}) = O(N \log N)$ , locality  $O(1)$ , and read efficiency  $O(1)$ .

### 3.3 An Optimal Instantiation for any Locality

As discussed in Section 1.1, the lower bound that we prove for schemes in the pad-and-split framework matches the tradeoff provided by the scheme of Demertzis and Papamanthou [DP17] (which is captured by our framework). Specifically, when setting the read efficiency of their scheme to  $O(1)$ , one obtains a statically-secure scheme with space  $O(N \log N / \log L)$ , locality  $L$ , and read efficiency  $O(1)$ . It should be noted that their scheme supports also non-constant read efficiency, but in that case it is not captured by our framework as it leaks additional information (in particular, the random choices made by the setup algorithm).

In what follows we describe their instantiation within our above-described template. Their scheme is obtained by splitting each list to sublists of lengths that are a power of the locality  $L$ . In our notation, we set  $m = \lfloor \log N / \log L \rfloor + 1 = \lfloor \log_L N \rfloor + 1$ ,  $\text{len}^{(j)} = L^{j-1}$ , and  $s^{(j)} = 2N / \text{len}^{(j)}$  for every  $j \in [m]$ . As for the splitting algorithm, a list of length  $L^{j-1} \leq n_i < L^j$  is padded to a length that is a multiple of  $L^{j-1}$ , and split into at most  $L$  sublists of length  $L^{j-1}$ . More formally,  $\text{SplitList}(N, n_i)$  outputs a vector of length  $m$ , where all the entries are zeros except for the entry in the position  $j = \lfloor \log_L(n_i) \rfloor + 1$ , which is set to the value  $\lceil n_i / L^{j-1} \rceil \in \{1, \dots, L\}$ .

This is indeed a valid  $\text{SplitList}$  algorithm, and its locality is  $L$ . Specifically, for each  $n_i$  and  $j$  it holds that  $\lceil n_i / L^{j-1} \rceil \cdot L^{j-1} \geq n_i$ , that is, each list is not longer than the sum of the lengths of its sublists. Moreover, for  $j = \lfloor \log_L(n_i) \rfloor + 1$  it also holds that  $\lceil n_i / L^{j-1} \rceil \leq L$  and  $\lceil n_i / L^{j-1} \rceil \cdot L^{j-1} < 2 \cdot n_i$ . This means that the locality is  $L$ , and that the padding at most doubles the length of the list. Therefore, it suffices to set  $s^{(j)} = 2N / \text{len}^{(j)}$ , and thus it holds that  $\sum_{j=1}^m \text{len}^{(j)} \cdot s^{(j)} = m \cdot 2N = O(N \cdot \log N / \log L)$ .

According to Theorem 3.3 and Theorem 3.4, the above splitting algorithm results in a searchable symmetric encryption schemes with space  $O(N \cdot \log N / \log L)$ , locality  $O(L)$ , and read efficiency  $O(1)$ . This yields the following corollaries:

**Corollary 3.7** ([DP17]). *Assuming the existence of any one-way function, for any  $L = L(N) > c$  (where  $c$  is an absolute constant) there exists a static  $\mathcal{L}_{\min}$ -secure symmetric searchable encryption scheme for databases of size  $N$  with the following parameters:*

- Space  $O(N \cdot \log N / \log L)$ .
- Locality  $L(N)$ .
- Read efficiency  $O(1)$ .
- Token size  $O(1)$ .

*Moreover, the scheme is adaptively  $\mathcal{L}_{\min}^{\text{adap}}$ -secure in the random-oracle model, when instantiating its building blocks appropriately.*

**Corollary 3.8.** *Assuming the existence of any one-way function, for any  $L = L(N) > c$  (where  $c$  is an absolute constant) there exists an adaptive  $\mathcal{L}_{\text{size}}^{\text{adap}}$ -secure symmetric searchable encryption scheme for databases of size  $N$  with the following parameters:*

- Space  $O(N \cdot \log N / \log L)$ .
- Locality  $L(N)$ .
- Read efficiency  $O(1)$ .
- Token size  $O(L(N))$ .

**Better efficiency for super-constant sub-polynomial locality.** For locality  $L(N)$  satisfying  $\omega(1) \leq L(N) \leq N^{o(1)}$  we can in fact instantiate our framework in a manner that reduces the expression  $\sum_{j=1}^m \text{len}^{(j)} s^{(j)}$  to  $(1 + o(1))(N \cdot \log N / \log L)$ . This matches our lower bound, which is shown to be  $(1 - o(1))(N \cdot \log N / \log L)$ , to within an additive lower-order term.

This is done as follows. Let  $\widehat{L} = \lfloor L / \log L \rfloor$ , and for a list of length  $n_i$  let  $j$  such that  $\widehat{L}^j \leq n_i < \widehat{L}^{j+1}$ . Represent  $n_i = a \cdot \widehat{L}^j + b \cdot \widehat{L}^{j-1} + c$ , where  $a \in \{1, \dots, \widehat{L} - 1\}$ ,  $b \in \{0, \dots, \widehat{L} - 1\}$ , and  $c \in \{0, \dots, \widehat{L}^{j-1} - 1\}$ . If  $a \geq \log L$ , then pad and split the list into at most  $\widehat{L}$  sublists of length  $\widehat{L}^j$ . Otherwise, pad and split the list into at most  $\widehat{L} \cdot \log L \leq L$  sublists of length  $\widehat{L}^{j-1}$ . This way, we never pad a list more than  $(1 + 1/\log L)$  times its length, since if  $a \geq \log L$  then the resulting padded length is upper bounded by

$$(a + 1) \cdot \widehat{L}^j = \left(1 + \frac{1}{a}\right) \cdot a \cdot \widehat{L}^j \leq \left(1 + \frac{1}{a}\right) \cdot n_i \leq \left(1 + \frac{1}{\log L}\right) \cdot n_i,$$

and if  $a < \log L$  then the resulting padded length is upper bounded by

$$\begin{aligned} a \cdot \widehat{L}^j + (b + 1) \cdot \widehat{L}^{j-1} &= \left(1 + \frac{1}{a \cdot \widehat{L} + b}\right) \cdot (a \cdot \widehat{L}^j + b \cdot \widehat{L}^{j-1}) \\ &\leq \left(1 + \frac{1}{a \cdot \widehat{L} + b}\right) \cdot n_i \\ &\leq \left(1 + \frac{1}{\log L}\right) n_i, \end{aligned}$$

where in the last inequality we used the fact that  $a \cdot \widehat{L} + b \geq \widehat{L} \geq \log L$  holds for sufficiently large  $L$ . Since we never pad a list more than  $(1 + 1/\log L)$  times its length, for any  $j$  we can set  $s^{(j)} = (1 + 1/\log L)N/\text{len}^{(j)}$ , and obtain

$$\sum_{j=1}^m \text{len}^{(j)} s^{(j)} = \left(1 + \frac{1}{\log L}\right) N \cdot \left(\left\lfloor \frac{\log N}{\log \widehat{L}} \right\rfloor + 1\right) = (1 + o(1))N \cdot \frac{\log N}{\log L},$$

where the last equality holds since  $\omega(1) \leq L \leq N^{o(1)}$ .

### 3.4 Our Lower Bound for Pad-and-Split Schemes

In this section we present our lower bound on the trade-off between the space and the locality of any pad-and-split scheme. Recall that each such a scheme is characterized by a **SplitList** algorithm that satisfies a modest validity requirement (recall Definition 3.1), and is associated with the following parameters (all of which may be functions of the size  $N$  of the database):

- The possible lengths  $\text{len}^{(1)}, \dots, \text{len}^{(m)}$  of sublists to which the **SplitList** algorithm splits the list associated with each keyword, as described in Section 3.1.
- Upper bounds  $s^{(1)}, \dots, s^{(m)}$  on the total number of sublists of lengths  $\text{len}^{(1)}, \dots, \text{len}^{(m)}$ , respectively, that are produced by the **SplitList** algorithm when processing an entire database.

Equipped with the above parameters, recall from Theorems 3.3 and 3.4 that the space usage of a pad-and-split scheme is  $O\left(\sum_{j=1}^m s^{(j)} \cdot \text{len}^{(j)}\right)$ , and the locality of such a scheme is  $O(L)$  where  $L = L(N)$  is the locality of its **SplitList** algorithm (i.e., each list is split into at most  $L$  sublists). Thus, proving a lower bound on the trade-off between the space and the locality of pad-and-split schemes translates to proving such a lower bound on the corresponding parameters of their underlying **SplitList** algorithm. We prove the following theorem from which Theorem 1.1 follows as an immediate corollary:

**Theorem 3.9.** *Let  $\text{SplitList}$  be a valid splitting algorithm with parameters  $\text{len}^{(1)}, \dots, \text{len}^{(m)}$  and  $s^{(1)}, \dots, s^{(m)}$ , and with locality  $L = L(N)$ . Then, for any  $0 < c < 1$  it holds that*

$$\sum_{j=1}^m \text{len}^{(j)} \cdot s^{(j)} \geq (1 - c) \cdot N \cdot \left( \frac{\log N}{\log L - \log c + C_1} - C_2 \right),$$

where  $C_1$  and  $C_2$  are small absolute constants.

In particular, by setting  $c = 1/2$  we obtain the lower bound  $\sum_{j=1}^m \text{len}^{(j)} \cdot s^{(j)} = \Omega(N \cdot \log N / \log L)$ , which implies Theorem 1.1. In addition, if  $\omega(1) \leq L(N) \leq N^{o(1)}$  then by setting  $c = 1/\log L$  we obtain the tighter lower bound  $\sum_{j=1}^m \text{len}^{(j)} \cdot s^{(j)} \geq (1 - o(1))N \cdot \log N / \log L$ .

The proof of Theorem 3.9 relies on two main observations that we formalize in Claims 3.10 and 3.11. First, in Claim 3.10 we prove that for every possible length  $1 \leq n_i \leq N$  of a list  $\text{DB}(w_i)$  in the database, when computing  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$  it holds that the contribution, in terms of space usage, of the sublists of lengths at least roughly  $n_i/L$  (where  $L$  is the locality) to which  $\text{DB}(w_i)$  is split must be linear in  $n_i$ . In other words, the  $\text{SplitList}$  algorithm must assign roughly  $n_i$  elements of the list  $\text{DB}(w_i)$  to sublists of length at least  $n_i/L$ .

For stating Claim 3.10 we introduce the following notation. Recall that  $\text{len}^{(1)}, \dots, \text{len}^{(m)}$  are the possible lengths of sublists to which the  $\text{SplitList}$  algorithm splits the list  $\text{DB}(w_i)$  associated with each keyword  $w_i$ . For any  $a < b \leq \infty$  we let

$$\text{lengths}[a, b] \stackrel{\text{def}}{=} \{j \in [m] \mid a < \text{len}^{(j)} \leq b\}.$$

That is, these are the indices of the possible lengths in the interval  $(a, b]$ .

**Claim 3.10.** *For any  $1 \leq n_i \leq N$  and  $0 < c < 1$  it holds that*

$$\sum_{j \in \text{lengths}[cn_i/L, \infty]} \text{len}^{(j)} \cdot x_i^{(j)} > (1 - c) \cdot n_i,$$

where  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$ .

**Proof.** Let  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$  and fix  $0 < c < 1$ . We show that there always exists a subset  $J \subseteq \text{lengths}[cn_i/L, \infty]$  of indices for which  $\sum_{j \in J} \text{len}^{(j)} \cdot x_i^{(j)} > (1 - c) \cdot n_i$ . This implies that

$$\sum_{j \in \text{lengths}[cn_i/L, \infty]} \text{len}^{(j)} \cdot x_i^{(j)} \geq \sum_{j \in J} \text{len}^{(j)} \cdot x_i^{(j)} > (1 - c) \cdot n_i$$

and thus the claim follows.

We construct the set  $J$  as follows, when initially setting  $J = \emptyset$ . Since  $\text{SplitList}$  is valid and has locality  $L$ , then it holds that  $\sum_{j=1}^m \text{len}^{(j)} \cdot x_i^{(j)} \geq n_i$  and  $\sum_{i=1}^m x_i^{(j)} \leq L$ . Therefore, there must exist some  $j_1 \in [m]$  for which  $x_i^{(j_1)} > 0$  and  $\text{len}^{(j_1)} \geq n_i/L$ . Otherwise, if for all  $j \in [m]$  such that  $x_i^{(j)} > 0$  it holds that  $\text{len}^{(j)} < n_i/L$ , we obtain the following contradictory inequality:

$$n_i \leq \sum_{j=1}^m \text{len}^{(j)} \cdot x_i^{(j)} < \sum_{j=1}^m \frac{n_i}{L} \cdot x_i^{(j)} \leq n_i.$$



We add this index  $j_1$  into the set  $J$  that we are constructing, pointing out that  $\text{len}^{(j_1)} \geq n_i/L > cn_i/L$  and thus  $j_1 \in \text{lengths}[cn_i/L, \infty]$  as required. Now, if  $\sum_{j \in J} \text{len}^{(j)} \cdot x_i^{(j)} > (1-c) \cdot n_i$  then we are done with constructing the set  $J$ . Otherwise, it holds that

$$\sum_{j \in [m] \setminus J} \text{len}^{(j)} \cdot x_i^{(j)} \geq n_i - (1-c) \cdot n_i \geq c \cdot n_i$$

In addition, it holds that  $\sum_{j \in [m] \setminus J} x_i^{(j)} < L$ , and from similar reasoning as before there must exist an index  $j_2 \in [m] \setminus J$  for which  $x_i^{(j_2)} > 0$  and  $\text{len}^{(j_2)} > c \cdot n_i/L$ . Therefore, we again conclude that  $j_2 \in \text{lengths}[cn_i/L, \infty]$ , and add the index  $j_2$  into the set  $J$ . We iteratively repeat this process, where in the  $k$ -th iteration either  $\sum_{j \in J} \text{len}^{(j)} \cdot x_i^{(j)} > (1-c) \cdot n_i$  or there exists  $j_k \in [m] \setminus J$  such that  $j_k \in \text{lengths}[cn_i/L, \infty]$ , and we add  $j_k$  to  $J$ . Since  $[m]$  is finite, this process will eventually terminate with either some iteration in which  $\sum_{j \in J} \text{len}^{(j)} \cdot x_i^{(j)} > (1-c) \cdot n_i$ , or with  $J = [m]$ . When  $J = [m]$ , then  $\sum_{j \in J} \text{len}^{(j)} \cdot x_i^{(j)} > (1-c)n_i$  must hold. To see that, recall that  $\sum_{j=1}^m \text{len}^{(j)} \cdot x_i^{(j)} \geq n_i$  and therefore if  $\sum_{j \in J} \text{len}^{(j)} \cdot x_i^{(j)} \leq (1-c)n_i$  we get a contradiction. In conclusion, we obtain

$$\sum_{j \in \text{lengths}[cn_i/L, \infty]} \text{len}^{(j)} \cdot x_i^{(j)} \geq \sum_{j \in J} \text{len}^{(j)} \cdot x_i^{(j)} > (1-c) \cdot n_i$$

as claimed. ■

Recall that a **SplitList** algorithm is parameterized by the possible lengths  $\text{len}^{(1)}, \dots, \text{len}^{(m)}$  of sublists, and also by upper bounds  $s^{(1)}, \dots, s^{(m)}$  on the total number of sublists of lengths  $\text{len}^{(1)}, \dots, \text{len}^{(m)}$ , respectively. Our next claim establishes lower bounds on the values  $s^{(1)}, \dots, s^{(m)}$ .

**Claim 3.11.** *For all  $1 \leq n_i \leq N$  and  $1 \leq j \leq m$ , it holds that*

$$s^{(j)} \geq x_i^{(j)} \cdot \left\lfloor \frac{N}{n_i} \right\rfloor$$

where  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$ .

**Proof.** Consider a database that consists of  $\lfloor N/n_i \rfloor$  lists of length  $n_i$ , and one additional list of length  $N - n_i \cdot \lfloor N/n_i \rfloor$  in case that  $N/n_i$  is not an integer. Recall that for **SplitList** to be valid, it must be that  $s^{(j)}$  upper bounds the number of sublists of length  $\text{len}^{(j)}$ . In our case, there are at least  $x_i^{(j)} \cdot \lfloor N/n_i \rfloor$  such sublists, so the claim follows. ■

Equipped with Claims 3.10 and 3.11, we are now ready to prove Theorem 3.9.

**Proof of Theorem 3.9.** We first prove the theorem under the assumption that each list  $\text{DB}(w_i)$  of length  $n_i$  is split into sublists of length at most  $3n_i$ . Formally, for every  $1 \leq n_i \leq N$  when computing  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$  it holds that  $x_i^{(j)} = 0$  for all  $j \in \text{lengths}[3n_i, \infty]$ . Claim 3.10 guarantees that for every  $1 \leq n_i \leq N$  it holds that

$$\sum_{j \in \text{lengths}[cn_i/L, 3n_i]} x_i^{(j)} \cdot \text{len}^{(j)} > (1-c)n_i, \tag{3.1}$$

where  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$ . Note that in the sum above we omitted lengths greater than  $3n_i$  due to our assumption. Claim 3.11 and Eq. (3.1) now imply that for every  $1 \leq n_i \leq N$  it holds that

$$\begin{aligned} \sum_{j \in \text{lengths}[cn_i/L, 3n_i]} s^{(j)} \cdot \text{len}^{(j)} &> \sum_{j \in \text{lengths}[cn_i/L, 3n_i]} x_i^{(j)} \cdot \left\lfloor \frac{N}{n_i} \right\rfloor \cdot \text{len}^{(j)} \\ &> \sum_{j \in \text{lengths}[cn_i/L, 3n_i]} x_i^{(j)} \cdot \frac{N - n_i}{n_i} \cdot \text{len}^{(j)} \\ &> (1 - c)(N - n_i). \end{aligned}$$

Let  $T(r) = \sum_{j \in \text{lengths}[c \cdot r/L, 3 \cdot r]} s^{(j)} \cdot \text{len}^{(j)}$  and  $z = \lceil 3L/c \rceil$ . When considering two ranges  $[c/L \cdot r_1, 3 \cdot r_1]$  and  $[c/L \cdot r_2, 3 \cdot r_2]$ , note that these two ranges do not intersect if  $r_2 \geq (3L/c) \cdot r_1$ . Therefore, we conclude that

$$\begin{aligned} \sum_{j=1}^m s^{(j)} \cdot \text{len}^{(j)} &\geq T(1) + T(z) + T(z^2) + \dots + T(z^{\lfloor \log N / \log z \rfloor - 1}) \\ &\geq (1 - c) \left( (N - 1) + (N - z) + \dots + (N - z^{\lfloor \log N / \log z \rfloor - 1}) \right) \quad (3.2) \\ &= (1 - c) \left( N \cdot \left\lfloor \frac{\log N}{\log z} \right\rfloor - \frac{z^{\lfloor \log N / \log z \rfloor - 1} - 1}{z - 1} \right) \\ &\geq (1 - c) \left( N \cdot \left\lfloor \frac{\log N}{\log z} \right\rfloor - \frac{N}{z - 1} \right) \\ &\geq (1 - c)N \cdot \left( \frac{\log N}{\log z} - 2 \right). \end{aligned}$$

For our choice of  $z = \lceil 3L/c \rceil$  it holds that  $\log z \leq \log L - \log c + 2$ , and thus the theorem follows.

We now remove the assumption that each list  $\text{DB}(w_i)$  of length  $n_i$  is split into sublists of length at most  $3n_i$  (i.e., the assumption that  $x_i^{(j)} = 0$  for all  $j \in \text{lengths}[3n_i, \infty]$ ). Recall that our above calculation followed from the fact that for each  $0 \leq k \leq \lfloor \log N / \log z \rfloor - 1$  it holds that  $T(z^k) \geq (1 - c)(N - z^k)$ , but without the above assumption this inequality is not guaranteed. Instead, we prove the following claim:

**Claim 3.12.** *For  $x \in \mathbb{R}$  denote  $(x)_+ \stackrel{\text{def}}{=} \max\{x, 0\}$ . Then, for every integer  $0 \leq k \leq \lfloor \log N / \log z \rfloor - 1$  there exists some integer  $k' \geq k$  such that*

$$\begin{aligned} T(z^k) + T(z^{k+1}) + \dots + T(z^{k'}) \\ \geq (1 - c) \left( (N - z^k)_+ + (N - z^{k+1})_+ + \dots + (N - z^{k'})_+ \right). \end{aligned}$$

Intuitively, Claim 3.12 guarantees that each  $T(z^k)$  is at least  $(1 - c)(N - z^k)_+$  in an amortized manner. Equipped with Claim 3.12, starting with  $k_0 = 0$  there exists some  $k_1 \geq 0$  for which

$$T(1) + T(z) + \dots + T(z^{k_1}) \geq (1 - c) \left( (N - 1)_+ + (N - z)_+ + \dots + (N - z^{k_1})_+ \right).$$

Next, as long as  $k_1 + 1 \leq \lfloor \log N / \log z \rfloor - 1$ , then again there exists some  $k_2 \geq k_1 + 1$  for which it holds that

$$T(z^{k_1+1}) + \dots + T(z^{k_2}) \geq (1 - c) \left( (N - z^{k_1+1})_+ + \dots + (N - z^{k_2})_+ \right).$$

Continuing this way, we define a sequence  $k_0 < k_1 < \dots < k_u$  such that  $k_u \geq \lceil \log N / \log z \rceil - 1$ . Summing all those up, we obtain that

$$\begin{aligned}
\sum_{j=1}^m s^{(j)} \cdot \text{len}^{(j)} &\geq T(1) + T(z) + T(z^2) + \dots + T(z^{k_u}) \\
&\geq (1-c) \left( (N-1)_+ + (N-z)_+ + \dots + (N-z^{k_u})_+ \right) \\
&\geq (1-c) \left( (N-1)_+ + (N-z)_+ + \dots + (N-z^{\lceil \log N / \log z \rceil - 1})_+ \right) \\
&= (1-c) \left( (N-1) + (N-z) + \dots + (N-z^{\lceil \log N / \log z \rceil - 1}) \right),
\end{aligned}$$

and following the same calculation as in Equation (3.2) above, we obtain that  $\sum_{j=1}^m s^{(j)} \cdot \text{len}^{(j)} \geq (1-c)N(\log N / \log z - 2)$  as before, and thus the theorem follows. We conclude the proof of the theorem by proving Claim 3.12.

**Proof of Claim 3.12.** Let  $n_i = z^k$  and  $(x_i^{(1)}, \dots, x_i^{(m)}) = \text{SplitList}(N, n_i)$ . If for every  $j \in \text{lengths}[3n_i, \infty]$  it holds that  $x_i^{(j)} = 0$ , then we already showed that  $T(z^k) \geq (1-c)(N - z^k)$ , and the claim holds for  $k' = k$ . Otherwise, let  $j \in [m]$  be any index such that  $\text{len}^{(j)} > 3n_i$  and  $x_i^{(j)} > 0$ , and let  $k' > k$  such that  $\text{len}^{(j)} \in \text{lengths}[c \cdot z^{k'}/L, 3 \cdot z^{k'}]$ . Claim 3.11 implies that

$$\begin{aligned}
T(z^{k'}) &\geq s^{(j)} \cdot \text{len}^{(j)} \\
&\geq x_i^{(j)} \cdot \frac{N - n_i}{n_i} \cdot 3 \cdot z^{k'-1} \\
&= 2 \cdot z^{k'-k-1} \cdot \frac{3}{2} \cdot (N - z^k).
\end{aligned}$$

Since  $k \leq \lceil \log N / \log z \rceil - 1$  and  $z \geq 4$ , it holds that  $z^k \leq N/z \leq N/3$ . Since  $k' > k$ , it holds that  $2 \cdot z^{k'-k-1} \geq (k' - k + 1)$ . Therefore, it holds that  $T(z^{k'}) \geq (k' - k + 1) \cdot N$ . We conclude that

$$\begin{aligned}
T(z^k) + T(z^{k+1}) + \dots + T(z^{k'}) &\geq T(z^{k'}) \\
&\geq (k' - k + 1) \cdot N \\
&\geq (1-c) \left( (N - z^k)_+ + (N - z^{k+1})_+ + \dots + (N - z^{k'})_+ \right)
\end{aligned}$$

as claimed. ■

## 4 The Statistical-Independence Framework: A Leveled Two-Choice Scheme

In this section we consider the statistical-independence framework introduced by Asharov et al. [ANS<sup>+</sup>16] for the design of symmetric searchable encryption schemes. As discussed in Section 1.2, within this framework we construct a scheme whose read efficiency when querying for a keyword  $w$  may depend on the length of the list  $\text{DB}(w)$  that is associated with  $w$ , and for any  $n \leq N$  we denote by  $r(N, n)$  the read efficiency when retrieving a list of length  $n$ .<sup>10</sup> We prove the following theorem:

<sup>10</sup>We emphasize that having the read efficiency depend on the length of the retrieved list does not hurt the security of SSE schemes, and our scheme still results in minimal leakage as required.

**Theorem 4.1.** *Assuming the existence of any one-way function, for any function  $f(N) = \omega(1)$  there exists an adaptive  $\mathcal{L}_{\text{size}}^{\text{adap}}$ -secure symmetric searchable encryption scheme for databases of size  $N$  in which no keyword is associated with more than  $N/\log^3 N$  identifiers, with the following parameters:*

- *Space  $O(N)$ .*
- *Locality  $O(1)$ .*
- *Read efficiency  $r(N, n) = f(N) \cdot \epsilon(n)^{-1} + O(\log \log \log N)$ , where  $n = N^{1-\epsilon(n)}$ .*
- *Token size  $O(1)$ .*

Comparing the performance of our new scheme with the lower bound of Asharov et al. in the statistical-independence framework, Theorem 4.1 matches their lower bound to within an additive  $O(\log \log \log N)$  factor in the read efficiency. Specifically, Asharov et al. proved the following lower bound for schemes in the statistical-independence framework (restated to consider read efficiency  $r(N, n)$  that may depend on the length  $n$  of each list, and to consider constant locality):

**Theorem 4.2** ([ANS<sup>+</sup>16]). *For any searchable symmetric encryption scheme in the statistical-independence framework with space  $O(N \log N)$ , locality  $O(1)$ , and read efficiency  $r(N, n)$ , there exists a function  $f(N) = \omega(1)$  such that  $r(N, n) = f(N) \cdot \epsilon(n)^{-1}$  for every  $1 \leq n \leq N/\log N$ , where  $n = N^{1-\epsilon(n)}$ .*

In the remainder of this section we first overview the statistical independence framework for the design of symmetric searchable encryption schemes (Section 4.1), and then present our new scheme within this framework (Section 4.2).

#### 4.1 The Statistical-Independence Framework

The statistical-independence framework of Asharov et al. [ANS<sup>+</sup>16] considers symmetric searchable encryption schemes that are characterized by a pair of algorithms, denoted **RangesGen** and **Allocation**, and consist of the following two phases:

- Given a database  $\text{DB} = \{\text{DB}(w_1), \dots, \text{DB}(w_{n_W})\}$  of size  $N$ , for each keyword  $w_i$  the scheme invokes the **RangesGen** algorithm on the length  $n_i$  of its corresponding list  $\text{DB}(w_i)$ , to obtain a set of *possible locations* in which the scheme may place the elements of the list  $\text{DB}(w_i)$ . This set consists of several intervals and we denote it by  $R_i = \{[a_1, b_1], \dots, [a_d, b_d]\} \leftarrow \text{RangesGen}(N, n_i)$ . Looking ahead, when supplied with a token corresponding to a keyword  $w_i$ , the server will return to the client all data stored in the possible locations of the list  $\text{DB}(w_i)$  (the server will not actually know in which of the possible locations the elements of the list are actually placed).
- Given the sets of possible locations  $R_1, \dots, R_{n_W}$  of the lists corresponding to all keywords  $w_1, \dots, w_{n_W}$ , respectively, the scheme invokes the **Allocation** algorithm on these sets (and on the respective lengths of the lists) to obtain the *actual locations* for the elements of all lists. We denote the actual locations as an array  $\text{map} \leftarrow \text{Allocation}((n_1, R_1), \dots, (n_{n_W}, R_{n_W}))$ , where each of its entries is either a pair  $(i, j)$  (representing that this entry is the actual location of the  $j$ th element from the list  $\text{DB}(w_i)$ ) or NULL (representing an empty entry).

A key property of this framework is that the **RangesGen** algorithm, which determines the set of possible locations for each list  $\text{DB}(w_i)$ , is applied separately and independently to the length of each list. Thus, the possible locations of each list are independent of the possible locations of all other lists (in contrast, the actual locations of the lists are naturally allowed to be correlated).

Asharov et al. referred to a pair (**RangesGen**, **Allocation**) of such algorithms as an allocation scheme, and showed that any such allocation scheme satisfying a natural correctness requirement can be used to construct a searchable symmetric encryption scheme. The correctness requirement asks that for any database, with all but a negligible probability, these algorithms produce an actual allocation **map** in which each element has exactly one actual placement (where the probability is taken over the internal coin tosses of the algorithms **RangesGen** and **Allocation**).

The resulting scheme of Asharov et al. inherits its space, locality and read efficiency from those of its underlying allocation scheme, defined as follows:

**Definition 4.3.** A pair (**RangesGen**, **Allocation**) of algorithms satisfying the above correctness requirement is an  $(s, d, r)$ -allocation scheme, for some functions  $s(\cdot)$ ,  $d(\cdot)$  and  $r(\cdot, \cdot)$ , if the following properties hold:

- **Space:** For any input  $(n_1, \dots, n_k)$ , the array  $\text{map} \leftarrow \text{Allocation}((n_1, R_1), \dots, (n_k, R_k))$ , where  $R_i = \{[a_1, b_1], \dots, [a_d, b_d]\} \leftarrow \text{RangesGen}(N, n_i)$  for every  $i \in [k]$ , is of size at most  $s(N)$ , where  $N = \sum_{i=1}^k n_i$ .
- **Locality:** For any input  $(N, n_i)$ , the algorithm **RangesGen** outputs at most  $d(N)$  ranges.
- **Read efficiency:** For any input  $(N, n_i)$  for the algorithm **RangesGen** it holds that:

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

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

Equipped with the above notation, Asharov et al. proved the following:

**Theorem 4.4** ([ANS<sup>+</sup>16]). *Given any  $(s, d, r)$ -allocation scheme and any one-way function, there exists an  $\mathcal{L}_{\text{size}}^{\text{adap}}$ -secure searchable symmetric encryption scheme for databases of size  $N$  with space  $O(s(N))$ , locality  $O(d(N))$ , and read efficiency  $O(r(N, \cdot))$ .*

**From allocation algorithms to SSE schemes.** We conclude our high-level description of the statistical-independence framework by briefly overviewing the generic transformation from allocation schemes to SSE scheme. The reader is referred to [ANS<sup>+</sup>16] for the complete formal description of this transformation.

In a nutshell, the client runs the **RangesGen** and the **Allocation** procedures as described above to obtain the actual allocation **map** of all elements. Then, the client encrypts each identifier from each list  $\text{DB}(w)$  in **map** with a key that is derived from the keyword  $w$  using a pseudorandom function. In addition, any unused entry in the array is filled with a uniform string of the appropriate length.

When issuing a query corresponding to a keyword  $w$ , the client asks the server to retrieve the encrypted content of all possible locations of the list  $\text{DB}(w)$ .<sup>11</sup> Since these locations are chosen independently at random, this does not reveal any additional information on the structure of the database except for the length of the queried list. The client then identifies the actual locations and decrypts the data by itself.

---

<sup>11</sup>The details here are quite subtle. The server obtains the pseudorandom key that was used to produce randomness for the relevant invocation of **RangesGen**. In addition, the server stores the lengths of the lists in an encrypted manner, and can only reveal the lengths of the already-queried lists. Knowing both the pseudorandom key and the list length allows the server to compute the possible locations of the list  $\text{DB}(w)$ .

## 4.2 Our Leveled Two-Choice Scheme

In this section we present our new allocation scheme from which Theorem 4.4 provides the searchable symmetric encryption schemes guaranteed by Theorem 4.1. Our scheme consists of the following three levels for storing the elements of any given database DB of size  $N$ :

- The first level, named the “two-choice array”, consists of the two-choice SSE scheme of Asharov et al. [ANS<sup>+</sup>16] but with *an exponentially improved read efficiency*. In this array, each list  $\text{DB}(w_i)$  can be stored in one out of two possible intervals of consecutive locations, in a manner that we describe below as part of our **Allocation** algorithm. However, unlike the scheme of Asharov et al. we do not store all of the  $N$  elements of the database in this array. Instead, the key observation underlying our new scheme is that when viewing this array as a collection of bins, then by allowing a few lists to “overflow” from this level to the second level (overall at most  $\widehat{N} = N/\log N$  elements will overflow with all but a negligible probability), we can reduce the maximal load of each bin from  $\tilde{O}(\log \log N)$  (as in [ANS<sup>+</sup>16]) to  $O(\log \log \log N)$ . This then translates into improving the read efficiency in this level from  $\tilde{O}(\log \log N)$  to  $O(\log \log \log N)$ .
- The second level, named the “cuckoo hashing level”, stores the vast majority of the elements that overflow from the first level. This level consists of roughly  $\log N$  cuckoo hashing tables (see Section 2.4), where the  $j$  hash table is designed to store at most  $\widehat{N}/2^j$  values each of which of size  $2^j$ . These values are the lists that overflow from the first level (the  $j$ th table will store overflowing lists of length roughly  $2^j$ ).
- The third level, named the “stash level”, consists of a cuckoo hashing stash for each of the second-level cuckoo hashing tables. The goal of introducing this level is to reduce the failure probably of cuckoo hashing from noticeable to negligible (see Section 2.4), which is essential for the security of the resulting SSE scheme.

This leveled structure of our allocation scheme, and thus of our SSE scheme, guarantees that the possible locations for a list  $\text{DB}(w)$  of length  $n$  are its two possible intervals in the two-choice array, its two locations in the  $j$ th cuckoo hashing table for  $j = \log n$ , and anywhere in the stash of the  $j$ th cuckoo hashing table. In what follows we formally describe our allocation scheme (see Algorithm 4.6), which we prove to have space  $O(N)$ , locality 5, and read efficiency  $\omega(1) \cdot \epsilon(n)^{-1} + O(\log \log \log N)$  when retrieving lists of length  $n = N^{1-\epsilon(n)}$ .

**Theorem 4.5.** *For any function  $f(N) = \omega(1)$ , Algorithm 4.6 describes an  $(O(N), 5, r(N, n))$ -allocation scheme for databases of size  $N$  in which no keyword is associated with more than  $N/\log^3 N$  identifiers, where  $r(N, n) = f(N) \cdot \epsilon(n)^{-1} + O(\log \log \log N)$  and  $n = N^{1-\epsilon(n)}$ .*

**Proof of Theorem 4.5.** We assume without loss of generality that  $f(N) = o(\log \log N)$  (since otherwise, we may take  $\tilde{f}(N) = \min(f(N), o(\log \log N))$  instead). For the two-choice part of the algorithm, we make use of the following theorem from [ANS<sup>+</sup>16].

**Theorem 4.7** ([ANS<sup>+</sup>16] Theorem 3.5 Part 1). *Let  $S \geq n_1$  be a bound on the maximal length, and let  $m$  be the number of bins. Consider the two-choice allocation algorithm. Then, 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  $\frac{8N}{m} + \log \log \frac{N}{S} + 2$ , where the level of an element is the load of its bin right after inserting the element (e.g., the first element that is interested to the bin has level 1).*

**ALGORITHM 4.6** (Our Allocation Scheme (RangesGen, Allocation)).

**Input:** A vector of integers  $(n_1, \dots, n_k)$  representing the lengths of the lists  $L_1, \dots, L_k$  in the database. We let  $N = \sum_{i=1}^k n_i$ ,  $\widehat{N} = N/\log N$ , and assume for concreteness that the  $n_i$ 's are powers of 2, and that  $n_1 \geq n_2 \geq \dots \geq n_k$ .

**Parameters:**

- A bound  $S = N/\log^3 N$  on the length of the longest list in the database.
- The number  $m = N/\log \log \log N$  of bins in the two-choice array (it is chosen as a power of 2 and such that  $m \geq n_1$ ).
- A bound  $\text{BinSize} = O(\log \log \log N)$  on the size of each bin in the two-choice array.
- Stash sizes  $s_0, \dots, s_t$  where  $t = \log S$  and  $s_j = f(N) \cdot \epsilon_j$  for every  $j \in [t]$ , where  $2^j = N^{1-\epsilon_j}$  and  $\omega(1) \leq f(N) \leq o(\log \log N)$  may be any pre-specified function.

**The memory layout.** The memory is partitioned into the following segments:

1.  $m$  bins  $B_0, \dots, B_{m-1}$ , each of size  $\text{BinSize}$ .
2. Hash tables  $H_0, \dots, H_t$ , where each hash table  $H_j$  is implemented as a cuckoo hash table for  $\widehat{N}/2^j$  data items of size  $2^j$  each with a stash of size  $s_j$ .

**The RangesGen algorithm.** On input  $N$  and  $n_i$ :

1. Uniformly sample  $\alpha_{i,1}, \alpha_{i,2} \leftarrow \{0, \dots, \frac{m}{n_i} - 1\}$ .  
Consider the two super bins  $\widetilde{B}_{\alpha_{i,1}} = (B_{n_i \cdot \alpha_{i,1} + j})_{j=0}^{n_i-1}$  and  $\widetilde{B}_{\alpha_{i,2}} = (B_{n_i \cdot \alpha_{i,2} + j})_{j=0}^{n_i-1}$ .
2. Sample two hash table locations  $\beta_{i,1}, \beta_{i,2}$  for the cuckoo hash table  $H_{\log n_i}$ .
3. The possible ranges  $R_i$  are (1) The above two super-bins; (2) The two cells  $\beta_{i,1}, \beta_{i,2}$  in the hashtable  $H_{\log n_i}$ ; (3) The stash of the table  $H_{\log n_i}$ .

**The Allocation algorithm.**

1. Initialize  $m$  empty bins  $B_0, \dots, B_{m-1}$ , and an empty set  $\text{LeftOvers}$ .
2. Initialize hash tables  $H_0, \dots, H_t$ , where each hash table  $H_j$  is implemented as a cuckoo hash table for  $\widehat{N}/2^j$  entries of size  $2^j$  with a stash of size  $s_j$ .
3. For every list  $L_i$  with size  $n_i$  and ranges  $R_i$ , reconstruct  $(\alpha_{i,1}, \alpha_{i,2})$  and  $(\beta_{i,1}, \beta_{i,2})$  from  $R_i$ , and place the list  $L_i$  as follows:
  - (a) Consider the two super bins  $\widetilde{B}_{\alpha_{i,1}} = (B_{n_i \cdot \alpha_{i,1} + j})_{j=0}^{n_i-1}$  and  $\widetilde{B}_{\alpha_{i,2}} = (B_{n_i \cdot \alpha_{i,2} + j})_{j=0}^{n_i-1}$ . Let  $\beta \in \{\alpha_{i,1}, \alpha_{i,2}\}$  be the index of the least loaded super bin among  $\widetilde{B}_{\alpha_{i,1}}$  and  $\widetilde{B}_{\alpha_{i,2}}$ , where the load of a super bin is defined as the sum of loads of the bins that constitutes that super bin. If the load of the bins in  $\widetilde{B}_\beta$  is  $\text{BinSize}$ , then add  $L_i$  to  $\text{LeftOvers}$ . Otherwise, place the list  $L_i$  in the super bin  $\widetilde{B}_\beta$ . That is, for every  $j = 0, \dots, n_i - 1$ , place the  $j$ th element of the list  $L_i$  in the bin  $B_{n_i \cdot \beta + j}$ .
  - (b) If the list was not placed, then insert  $L_i$  into the cuckoo hash table  $H_{\log n_i}$  using the locations  $\beta_{i,1}$  and  $\beta_{i,2}$ . Note that the list might be placed in the stash. If the insertion fails, then output  $\perp$  and abort.

In Algorithm 4.6, we set  $S = N/\log^3 N$ ,  $m = N/\log \log \log N$ , and  $\text{BinSize} = O(\log \log \log N)$ . Therefore, with an overwhelming probability there are at most  $\hat{N} = N/\log N$  overflowing elements, and in this case, we place at most  $\hat{N}$  elements in the cuckoo hashing tables with the stashes.

Now we analyze the placement of the elements in the hash tables, assuming that the number of elements in `LeftOvers` is at most  $\hat{N}$ . For each  $0 \leq j \leq t$ , we set the stash size  $s_j = f(N) \cdot \epsilon_j^{-1}$  where  $\epsilon_j$  is chosen such that  $2^j = N^{1-\epsilon_j}$ . We obtain that the algorithm fails to insert the lists into the cuckoo hash table  $H_j$  with its stash with probability at most  $O((\hat{N}/2^j)^{-s_j/2})$  (see Section 2.4). Note that  $N^{\epsilon_j} \geq \log^3 N$ , so it holds that

$$\begin{aligned} (\hat{N}/2^j)^{-s_j/2} &= (N^{\epsilon_j}/\log N)^{-s_j/2} \\ &\leq (N^{\frac{2}{3}\epsilon_j})^{-s_j/2} \\ &= N^{-f(N)/3}. \end{aligned}$$

Thus, the insertion of overflowing elements fails with a negligible probability, and we conclude that Algorithm 4.6 fulfills the correctness requirement. Regarding read efficiency, the overhead of the 2-choice is  $O(\log \log \log N)$ , the overhead of the cuckoo hash table is 2, and the overhead of the stash is  $f(N) \cdot \epsilon(n)^{-1}$ , where  $n = N^{1-\epsilon(n)}$ , so in total we get an overhead of  $f(N) \cdot \epsilon_i^{-1} + O(\log \log \log N)$  as claimed. Locality of 5 easily follows from the description of `SplitList`. Regarding the space overhead, the bins require space of  $m \cdot \text{BinSize} = O(N)$ , each cuckoo hash table with stash requires space of  $O(\hat{N}) = O(N/\log N)$ , and there are less than  $\log N$  tables. So in total, the space overhead is  $O(N)$ . ■

## References

- [ADW14] M. Aumüller, M. Dietzfelbinger, and P. Woelfel. Explicit and efficient hash families suffice for cuckoo hashing with a stash. *Algorithmica*, 70(3):428–456, 2014.
- [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.
- [ANS<sup>+</sup>16] G. Asharov, M. Naor, G. Segev, and I. Shahaf. Searchable symmetric encryption: Optimal locality in linear space via two-dimensional balanced allocations. In *Proceedings of the 48th Annual ACM Symposium on Theory of Computing*, pages 1101–1114, 2016.
- [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.
- [DP17] I. Demertzis and C. Papamanthou. Fast searchable encryption with tunable locality. In *Proceedings of the 2017 ACM Special Interest Group on Management of Data (SIGMOD) Conference*, pages 1053–1067, 2017.
- [DPP17] I. Demertzis, D. Papadopoulos, and C. Papamanthou. Searchable encryption with optimal locality: Achieving sublogarithmic read efficiency. Cryptology ePrint Archive, Report 2017/749, 2017.
- [Goh03] E. Goh. Secure indexes. Cryptology ePrint Archive, Report 2003/216, 2003.
- [Gol04] O. Goldreich. Foundations of Cryptography – Volume 2: Basic Applications. Cambridge University Press, 2004.
- [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.
- [KMW09] A. Kirsch, M. Mitzenmacher, and U. Wieder. More robust hashing: Cuckoo hashing with a stash. *SIAM Journal on Computing*, 39(4):1543–1561, 2009.
- [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.
- [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.
- [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.
- [SPS14] E. Stefanov, C. Papamanthou, and E. Shi. Practical dynamic searchable encryption with small leakage. In *Proceedings of the 21st Annual Network and Distributed System Security Symposium*, 2014.
- [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.