

Dynamic Searchable Encryption in Very-Large Databases: Data Structures and Implementation

David Cash*, Joseph Jaeger*, Stanislaw Jarecki†, Charanjit Jutla‡,
Hugo Krawczyk‡, Marcel-Cătălin Roşu‡, and Michael Steiner‡

*Rutgers University

†University of California, Irvine

‡IBM Research

Abstract—We design and implement dynamic symmetric searchable encryption schemes that efficiently and privately search server-held encrypted databases with tens of billions of record-keyword pairs. Our basic theoretical construction supports single-keyword searches and offers asymptotically optimal server index size, fully parallel searching, and minimal leakage. Our implementation effort brought to the fore several factors ignored by earlier coarse-grained theoretical performance analyses, including low-level space utilization, I/O parallelism and goodput. We accordingly introduce several optimizations to our theoretically optimal construction that model the prototype’s characteristics designed to overcome these factors. All of our schemes and optimizations are proven secure and the information leaked to the untrusted server is precisely quantified. We evaluate the performance of our prototype using two very large datasets: a synthesized census database with 100 million records and hundreds of keywords per record and a multi-million webpage collection that includes Wikipedia as a subset. Moreover, we report on an implementation that uses the dynamic SSE schemes developed here as the basis for supporting recent SSE advances, including complex search queries (e.g., Boolean queries) and richer operational settings (e.g., query delegation), in the above terabyte-scale databases.

I. INTRODUCTION

BACKGROUND. Searchable symmetric encryption (SSE) allows one to store data at an untrusted server and later search the data for records (or documents) matching a given keyword while maintaining privacy. Many recent works [3]–[5], [7], [9], [14], [15], [17], [19], [21] studied SSE and provided solutions with varying trade-offs between security, efficiency, and the ability to securely update the data after it has been encrypted and uploaded. These constructions aim at practical efficiency, in contrast to generic cryptographic tools like homomorphic encryption or multiparty computation which are highly secure but not likely to be efficient in practice.

Large data sizes motivate storage outsourcing, so to be useful an SSE scheme must scale well. Existing SSE schemes employ only symmetric cryptography operations and standard

data structures and thus show potential for practical efficiency, but obstacles remain. While most constructions have theoretically optimal search times that scale only with the number of documents matching the query, the performance of their implementations on large datasets is less clear. Factors like I/O latency, storage utilization, and the variance of real-world dataset distributions degrade the practical performance of theoretically efficient SSE schemes. One critical source of inefficiency in practice (often ignored in theory) is a complete lack of locality and parallelism: To execute a search, most prior SSE schemes *sequentially* read each result from storage at a pseudorandom position, and the only known way to avoid this *while maintaining privacy* involves padding the server index to a prohibitively large size.

CONTRIBUTIONS. We give the first SSE implementation that can encrypt and search on datasets with tens of billions of record/keyword pairs. To design our scheme, we start with a new, simple, theoretical SSE construction that uses a generic dictionary structure to already achieve an asymptotic improvement over prior SSE schemes, giving optimal leakage, server size, search computation, and parallelism in search. This starting point can be seen as a generalization and simplification of the more ad-hoc techniques of [3]. We show how to make the scheme *dynamic*, meaning that the data can be changed after encryption: Our scheme can easily support additions to the data, as well as deletions via revocation lists.

Because the scheme uses a generic dictionary that itself has no security properties, it allows for several extensions and modifications with only small changes to the security proofs. In particular, our implementation effort showed that disk I/O utilization remained a bottleneck which prevented scaling; so we extend our basic construction to improve locality and throughput. These extensions preserve privacy with slightly different leakages that we analyze with formal security proofs. Below we describe the techniques behind results in more detail, starting with the new theoretical scheme that we extend later, and then compare our results to prior work.

BASIC CONSTRUCTION. Our scheme is very simple (see Figure 2): It associates with each record/keyword pair a pseudorandom label, and then for each pair stores the encrypted record identifier with that label in a generic dictionary data structure. We derive the labels so that the client, on input a keyword to query, can compute a keyword-specific short key allowing the server to search by first recomputing the labels, then retrieving the encrypted identifiers from the dic-

tionary, and finally decrypting the matching encrypted record identifiers. The only information leaked to the server by the encrypted index (other than the indexes of records matching a query) is the number of items in the dictionary, i.e. the number of record/keyword pairs in the data. This scheme is easy to implement correctly (and with parallel searching) because we make no security demands on the dictionary thus allowing instantiations as applications demand.

EXTENSIONS FOR EXTERNAL STORAGE. To compute the results of a keyword search with r matches, our basic scheme requires r retrievals from the dictionary for pseudorandom labels. Assuming $O(1)$ cost of a dictionary retrieval, this is asymptotically optimal. However, in implementations this will be far from optimal when the dictionary is stored in external memory (i.e., a block device like a HDD), because each random-looking retrieval will generate a disk read. This is in contrast to a plaintext system which could store all of the matches in a single contiguous area of memory.

In view of this reality we extend our scheme to use external storage more carefully while maintaining privacy. We first show how to securely “pack” related results together via a padding strategy to reduce the number of dictionary retrievals.

We found that even this modification was too slow for the datasets we targeted, and in particular we noticed that real data-sets exhibit extreme variability in the number of matches for a keyword: There were typically many keywords matching very few documents, then some keywords matching a significant fraction of the entire database. Our padding strategy thus becomes unsatisfactory because the (many) keywords matching only a few results create a lot of padding, and the searches that return a large number of results still trigger a large number of dictionary retrievals.

To address this we introduce further modifications that replace dictionary reads with array reads when processing large numbers of results. These modifications result in a slightly different, but intuitively acceptable (and perhaps even better) leakage profile that we discuss below.

EXTENSION FOR UPDATES. We observe that our scheme easily extends to allow for additions to the data after it has been uploaded. We only have to arrange that the client can compute the labels for the new data to be added, which it sends to the server for to be added to the dictionary. This requires either client state or communication proportional to the total number of keywords ever added or deleted. To support deletions we maintain a (pseudorandom) revocation list at the server that allows to filter out results that should be deleted; To actually reclaim space we must re-encrypt periodically.

OTHER APPLICATIONS. Recent constructions of SSE supporting more complex queries [3] and multi-client settings [13] use SSE as a black-box. Thus our data structures and associated operations (including support for dynamic databases) are readily available to support terabyte-scale databases in these much richer/complex encrypted-search settings (see end of Section II).

IMPLEMENTATION. Our implementation remains efficient on *two orders of magnitude larger* datasets than the most scalable

previous work [3], resulting in the first implementation of SSE on terabyte-scale databases containing tens of billions of indexed record/keyword pairs. We report on our prototype design and experimental results in Section V.

COMPARISON TO PRIOR WORK. In Figure 1 we compare our basic theoretical scheme to prior work. The basic scheme Π_{bas} generalizes and greatly simplifies an approach implicit in [3], which complicated the analysis by demanding security properties of the underlying data structures.

For a database with N record/keyword pairs, our basic scheme Π_{bas} produces an encrypted index of optimal size $O(N)$, leaks only the size N and the matching record id’s, and processes a search with r results in optimal $O(r)$ time, assuming $O(1)$ -cost for dictionary retrievals. Searching is trivial to parallelize with any number of processors.

Most prior schemes leak additional information, like the number of unique keywords, the size of the largest number of matches for a keyword, and so on. Some of these works also pad their encrypted indexes to be (worst-case) quadratic in their input size, which is totally impractical for large datasets. A notable issue with most prior work was a difficulty with parallelism: Other than [3], parallel searching was only achieved by two works that needed quadratic padding. Works like [7] required walking through an encrypted linked list and were not parallelizable at all. See the “Ind Leak”, “Index Size”, and “Search Time” columns in Figure 1.

The only prior dynamic schemes either had an impractically large index [14] or leaked the *structure* of the added documents [15], meaning that the server learned, say, the pattern of which keywords appear in which documents as they are added, which is a severe form of leakage compared to the usual SSE leakage of facts like the total database size. Our dynamic extension maintains the optimal index size and only leaks basic size information (and not document structure, as in [15]). Unlike prior dynamic schemes, ours does not reclaim space after each deletion - rather, we envision applications where deletions are relatively rare or, more generally, where a periodic complete re-encryption of the data is performed (re-encryption may be desirable to mitigate the leakage from updates with any dynamic SSE scheme).

MORE ON RELATED WORK. The notion of SSE we consider has its origins in work by Song, Wagner, and Perrig [19]. Several schemes since have improved upon the security and efficiency offered by the original schemes. The most similar to our construction is that of Chase and Kamara [5], and Cash et al [3]. Chase and Kamara also uses a dictionary, but in a different way and with an impractical level of padding for large datasets. Cash et al implements a scheme similar to our basic construction, but do not address updates nor, as we show in Section V-E, does their approach achieve the same level of practical scalability.

There is also a related line of work on searchable public-key encryption starting with [2], all of which do not scale due to linear-time searching. The version of SSE we deal with inherently leaks the identifiers of documents that match a query, as well as when a query is repeated. It is possible to hide even this information using private information retrieval [6] or oblivious RAM [10]. Several recent improvements to oblivious

Scheme	Security	Ind Leak	Dyn.?	Dyn Leak	Index Size	Search Time/Comm	Dyn. Comm
CGKO'06-1 [7]	NonAd	m, N	No	—	$O(N + m)$	$O(r), O(1)$	—
CGKO'06-2 [7]	Ad	Mn	No	—	$O(Mn)$	$O(r), O(r)$	—
CK'10 [5]	Ad	m, n, M	No	—	$O(Mn)$	$O(r), O(r)$	—
LSDHJ'10 [21]	Ad	m, n	Yes	no proof	$O(mn)$	$O(m), O(1)$	$O(W_{id})$
KO'12 [17]	Ad(UC)	n, M	No	—	$O(Mn)$	$O(r), O(1)$	—
KPR'12 [15]	Ad ^{ro}	m, N	Yes	EP(W_{id})	$O(N + m)$	$O(r), O(1)$	$O(W_{id})$
KP'13 [14]	Ad ^{ro}	m, n	Yes	minimal	$O(mn)$	$O((r \log n)/p), O(1)$	$O(W_{id} + m \log n)$
Basic (Π_{bas} here)	NonAd, Ad ^{ro}	N	No	—	$O(N)$	$O(r/p), O(1)$	—
Basic Adp (Π_{bas}^{ro} here)	Ad	N	No	—	$O(N)$	$O(r/p), O(r)$	—
Basic Dyn ($\Pi_{bas}^{dyn}, \Pi_{bas}^{dyn,ro}$ here)	NonAd, Ad ^{ro}	N	Yes	minimal	$O(N)$	$O((r + d_w)/p), O(1)$	$O(W_{id} + m \log n)$

Fig. 1. Comparison of some SSE schemes. Many leakages can be replaced by upper bounds and some search times assume interaction when the original paper was non-interactive. Legend: In security, “Ad” means adaptive security, Ad^{ro} means adaptive security in the random oracle model, and NonAd means non-adaptive security. Ind Leakage is leakage from encrypted database only. Search comm. is the size of the message sent from client ($O(r)$ from the server is inherent.) ro means random oracle model, $n = \#$ documents, $N = \sum_w |\text{DB}(w)|$, $m = |W|$, $M = \max_w |\text{DB}(w)|$, $r = |\text{DB}(w)|$ for the query w , $p = \#$ processors, $|W_{id}| = \#$ keyword changes in an update, EP(W_{id}) = structural equality pattern of changed keywords (see discussion at the end of Section IV), $d_w =$ the number of times the searched-for keyword has been added/deleted.

RAM move it further towards practicality [11], [20], but it is far from clear that they are competitive with SSE schemes at scale as one must implement the plaintext searching as an abstract RAM program and then run this program underneath the ORAM scheme without leaking information, say, via timing.

ORGANIZATION. Preliminary definitions are given in Section II. Our non-dynamic (i.e., static) SSE constructions are given in Section III, and the dynamic extensions are given in Section IV. Finally we report on our implementation in Section V.

II. DEFINITIONS AND TOOLS

The security parameter is denoted λ . We will use the standard notions of variable-input-length PRFs and symmetric encryption schemes (c.f. [16]). For these primitives we make the simplifying assumption that their keys are always in $\{0, 1\}^\lambda$, and that key generation for the encryption scheme picks a random key. Some of our constructions will be analyzed in the random oracle model [1], with the random oracle denoted H .

Our constructions will use a symmetric encryption scheme that is assumed to have *pseudorandom ciphertexts under chosen-plaintext attack*, which we call *RCPA security*. Formally, this says that an adversary cannot distinguish an oracle returning encryptions of chosen messages from one returning a random string with length equal to a fresh ciphertext for the chosen message.¹

SSE SCHEMES. We follow the formalization of Curtmola et al. [7] with some modifications discussed below. A database $\text{DB} = (\text{id}_i, W_i)_{i=1}^d$ is a list of identifier/keyword-set pairs where $\text{id}_i \in \{0, 1\}^\lambda$ and $W_i \subseteq \{0, 1\}^*$. When the DB under consideration is clear, we will write $W = \bigcup_{i=1}^d W_i$. For a keyword $w \in W$, we write $\text{DB}(w)$ for $\{\text{id}_i : w \in W_i\}$. We will always use $m = |W|$ and $N = \sum_{w \in W} |\text{DB}(w)|$ to mean the number of keywords and the total number of keyword/document matches in DB.

¹Our constructions can be proved secure assuming a type of key anonymity property, but RCPA is simpler and is anyway achieved by many efficient constructions.

A *dynamic searchable symmetric encryption (SSE) scheme* Π consists of an algorithm Setup and protocols Search and Update between the client and server, all fitting the syntax below. We assume that the server is deterministic, and that the client may hold some state between queries. A *static* SSE scheme is exactly the same, but with no Update protocol. Setup takes as input a database DB, and outputs a secret key K along with an encrypted database EDB. The search protocol is between a *client* and *server*, where the client takes as input the secret key K and a query $w \in \{0, 1\}^*$ and the server takes as input EDB and the server outputs a set of identifiers and the client has no output. In the Update protocol the client takes as input a key K , an operation $\text{op} \in \{\text{add}, \text{del}, \text{edit}^+, \text{edit}^-\}$, a file identifier id , and a set W_{id} of keywords. These inputs represent the actions of adding a new file with identifier id containing keywords W_{id} , deleting the file with identifier id , or add/removing the keywords in W_{id} from an existing file. At the end of the Update, the server outputs an updated encrypted database, and the client has no output.

We say that an SSE scheme is *correct* if the search protocol returns the (current) correct results for the keyword being searched (i.e., $\text{DB}(w)$), except with negligible probability. (More precisely, this should hold for every database DB, after a polynomially unbounded number of updates. We defer details to the full version.) To simplify the formalism we ignore the case where the client attempts to add a file with an existing identifier or delete/edit with an identifier that is not present in DB. Our protocols (and implementations) can handle these cases cleanly.

DISCUSSION. For simplicity our formalization of SSE does not model the storage of the actual document payloads. The SSE literature varies on its treatment of this issue, but in all cases one can augment the schemes to store the documents with no additional leakage beyond the length of the payloads. Compared to others we model also modifications of documents ($\text{edit}^+, \text{edit}^-$) in addition to add and delete of complete documents (add, del) as this can lead to more efficient protocols with reduced leakage.

The correctness definition for SSE requires the server to learn the ids of the results. One could define correctness to require the *client* to learn the ids instead. The two approaches are essentially equivalent assuming that encrypted documents

are of fixed length.

SECURITY. Security [5], [7], [15] follows the real/ideal simulation paradigm and is parametrized by a *leakage function* \mathcal{L} that describes what a secure protocol is allowed to leak. Formally, \mathcal{L} is a stateful party in an ideal security game, which is defined below.

Definition 1: Let $\Pi = (\text{Setup}, \text{Search}, \text{Update})$ be a dynamic SSE scheme and let \mathcal{L} be a leakage function. For algorithms A and S , define games $\mathbf{Real}_A^\Pi(\lambda)$ and $\mathbf{Ideal}_{A,S}^\Pi(\lambda)$ as follows:

$\mathbf{Real}_A^\Pi(\lambda)$: $A(1^\lambda)$ chooses DB. The game then runs $(K, \text{EDB}) \leftarrow \text{Setup}(\text{DB})$ and gives EDB to A . Then A repeatedly requests to engage in the Search or Update protocols, where A picks a client input in . To respond, the game runs the Search or Update protocol with client input (K, in) and server input EDB and gives the transcript to A (the server is deterministic so this constitutes its entire view). Eventually A returns a bit that the game uses as its own output.

$\mathbf{Ideal}_{A,S}^\Pi(\lambda)$: $A(1^\lambda)$ chooses DB. The game runs $\text{EDB} \leftarrow S(\mathcal{L}(\text{DB}))$ and gives EDB to A . Then A repeatedly requests to engage in the Search or Update protocols, where A picks a client input in . To respond, the game gives the output of $\mathcal{L}(\text{in})$ to S , which outputs a simulated transcript that is given to A . Eventually A returns a bit that the game uses as its own output.

Π is \mathcal{L} -secure against adaptive attacks if for all adversaries A there exists an algorithm S such that

$$\Pr[\mathbf{Real}_A^\Pi(\lambda) = 1] - \Pr[\mathbf{Ideal}_{A,S}^\Pi(\lambda) = 1] \leq \text{neg}(\lambda).$$

We define \mathcal{L} -security against non-adaptive attacks in the same way, except that in both games A must choose all of its queries at the start, \mathcal{L} takes them all as input, and S uses the output of \mathcal{L} to generate EDB and the transcripts at the same time. We also obtain adaptive and non-adaptive security definitions for static SSE schemes by disallowing adversary queries for Update.

DATA STRUCTURES. Our constructions will employ the standard data structures of lists, arrays, and dictionaries. We formalize a dictionary data type in detail because its syntax is relevant to our security analyses. Below, when we say *label*, *data*, or *data structure*, we mean *bitstring* and will treat them as such in the analysis.

An *dictionary implementation* D consists of four algorithms Create, Get, Insert, Remove. Create takes a list of label-data pairs $(\ell_i, d_i)_{i=1}^m$, where each label is unique, and outputs the data structure γ . On input γ and a label ℓ , Get(γ, ℓ) returns the data item with that label. On input γ and (ℓ, d) , Insert($\gamma, (\ell, d)$), outputs an updated data structure, that should contain the new pair. On input γ and ℓ , Remove(γ, ℓ) outputs an updated data structure with the pair removed.

We define *correctness* in the obvious way, i.e., the output of Get is always the data with the (unique) label it is given as input, and that it returns \perp when no data with the label is present.

We say that a dictionary implementation is *history-independent* if for all lists L the distribution of Create(L) depends only on the members of L and not their order in the list. The Create algorithm may be randomized or deterministic and satisfy history-independence. This simplest way to achieve it is to sort L first, but for large lists this will be infeasible.

We note that we only need the ability to remove data in some limited uses of dictionaries. In all settings where we need a very large dictionary, we can use an add-only version of the data structure.

EXTENSIONS AND GENERALIZATION. Two works [3], [13] showed that data structures for single-keyword SSE can be generalized to work for more complex SSE functionalities and models. Specifically, [3] shows how to extend SSE data structures to perform boolean queries on encrypted data (via the OXT protocol), and [13] further extends this functionality to more complex multi-user SSE settings. As a result, all the data structures in this paper, including their associated operations, their extreme efficiency and dynamic (updatable) capabilities, can be readily used to support these richer functional settings. All that is needed is to extend the data stored in these data structures from simple document identifiers (in the basic SSE case) to other associated data such as an encrypted key in the case of multi-client SSE (a key used by clients to decrypt documents) or protocol-specific values (such as the ‘ y ’ value in the OXT protocol from [3]). As a consequence, our advancement on the practicality and scale of SSE data structures immediately translates into the ability to support very large and dynamic databases even for functionalities as involved as full boolean SSE search in single- and multi-client SSE settings. We provide concrete evidence of this practical impact in Section V-E where we report performance numbers on query execution in these complex settings.

III. STATIC CONSTRUCTIONS

Let $D = (\text{Create}, \text{Get}, \text{Insert}, \text{Remove})$ be a dictionary implementation, F be a variable-input-length PRF, and $\Sigma = (\text{Enc}, \text{Dec})$ be a symmetric-key encryption scheme.

BASIC CONSTRUCTION. In Figure 2 we give our first and simplest construction, denoted Π_{bas} . To build the encrypted database, Setup(DB) chooses a key K and uses it to derive per-keyword keys for a PRF (to derive pseudorandom labels) and for encryption (to encrypt the identifiers). Then for each keyword w , it iterates over the identifiers in $\text{DB}(w)$. For each identifier, it computes a pseudorandom label by applying the PRF to a counter, encrypts the identifier, and adds the label/ciphertext pair to a list L . After all of the results have been processed it builds the dictionary γ from L , which becomes the server’s index. It is important that L is sorted by the labels before being loaded into the dictionary, or that the dictionary satisfies history independence - Without one of these, the scheme will leak information about the order in which the input was processed.

To search for keyword w , the client re-derives the keys for w and sends them to the server, who re-computes the labels and retrieves and decrypts the results.

LEAKAGE FUNCTION. The leakage function \mathcal{L} for our first construction responds to an initial startup query, and to search

Setup(DB)

1. $K \xleftarrow{\$} \{0, 1\}^\lambda$ allocate list L
2. For each $w \in W$:
 - $K_1 \leftarrow F(K, 1\|w)$, $K_2 \leftarrow F(K, 2\|w)$
 - Initialize counter $c \leftarrow 0$
 - For each $\text{id} \in \text{DB}(w)$:
 - $\ell \leftarrow F(K_1, c)$; $d \leftarrow \text{Enc}(K_2, \text{id})$; $c++$
 - Add (ℓ, d) to the list L (in lex order)
 - Set $\gamma \leftarrow \text{Create}(L)$
3. Output the client key K and $\text{EDB} = \gamma$.

Search

- Client:* On input (K, w) ,
- $K_1 \leftarrow F(K, 1\|w)$, $K_2 \leftarrow F(K, 2\|w)$
 - Send (K_1, K_2) to the server.
- Server:* For $c = 0$ until Get returns \perp ,
- $d \leftarrow \text{Get}(\gamma, F(K_1, c))$; $m \leftarrow \text{Dec}(K_2, d)$
 - Parse and output id in each m

Fig. 2. Scheme Π_{bas} .

queries, where its behavior is defined as follows. We describe the interactive stateful leakage function for the adaptive definitions; The non-adaptive leakage function is the obvious version that iterates over the queries with the adaptive leakage function. On initial input DB , \mathcal{L} outputs $N = \sum_{w \in W} |\text{DB}(w)|$, saves DB and an empty list Q_{srch} as state. Then, for a search input w , \mathcal{L} increments i , adds (i, w) to Q_{srch} and outputs $\text{DB}(w)$ and a set $\text{sp}(w, Q_{\text{srch}})$, called the *search pattern for w* , defined by

$$\text{sp}(w, Q_{\text{srch}}) = \{j : (j, w) \in Q_{\text{srch}}\}.$$

The search pattern indicates which other queries were also for the keyword w , and represents the fact that our scheme will send the same message when a search is repeated.

We deal with non-adaptive \mathcal{L} -security first.

Theorem 2: Π_{bas} is correct and \mathcal{L} -secure against non-adaptive attacks if F is a secure PRF and (Enc, Dec) is RCPA-secure.

Proof sketch: Without loss of generality we assume an adversary never repeats a query, as this obviously will not help because our search protocol is deterministic. Correctness holds because collisions amongst the labels produced by F occur with negligible probability in a random function and hence also with the PRF F . To prove non-adaptive security we must give a simulator that takes as input the leakage output setting up DB (which is N) as well as the leakages on each of the queries (which are just the sets $\text{DB}(w_i)$ of identifiers matching the i -th query, assuming all queries are unique). We need to produce the view of the server, which is the EDB data structure along with the messages from the client from each search.

The simulator iterates over the queries, choosing keys $K_1^i, K_2^i \xleftarrow{\$} \{0, 1\}^\lambda$ for the i -th query, and then for each $\text{id} \in \text{DB}(w_i)$ it computes ℓ and d as specified in the real Setup (using K_1^i and K_2^i as K_1 and K_2), adding each of the pairs to a list L . Then it adds random label/data pairs to L (still maintained in lexicographic order) until it has N total elements, and creates a dictionary $\gamma = \text{Create}(L)$. The simulator outputs $\text{EDB} = \gamma$ and the transcript for the i query is (K_1^i, K_2^i) .

A simple hybrid argument shows that the simulator's output is indistinguishable from the real server view. The first hybrid shows that selecting each K_1, K_2 at random is indistinguishable from deriving them from K , by the PRF security of F . The next hybrid shows that the labels and data ciphertexts for unqueried keywords are pseudorandom. The only subtlety is to verify that, because the list is maintained in lexicographic order, the relationship between the unopened result sets is not needed. \square

It is easy to prove secure (via a reduction to the above theorem) a version of Π_{bas} that does not maintain a sorted list L , as long as the dictionary is history-independent.

In the random oracle model we can achieve adaptive security for the same \mathcal{L} if F is replaced with the random oracle H so $F(K, x) := H(K\|x)$, and the encryption algorithm Enc , on inputs $K, m \in \{0, 1\}^\lambda$, chooses a random $r \in \{0, 1\}^\lambda$ and outputs $(r, H(K\|r) \oplus m)$. We denote this variant $\Pi_{\text{bas}}^{\text{ro}}$.

Theorem 3: $\Pi_{\text{bas}}^{\text{ro}}$ is \mathcal{L} -secure against adaptive attacks in the random oracle model.

Proof sketch: This theorem is proved in a similar way to the previous one, except that the simulator programs the random oracle in response to adaptive queries to open the labeled ciphertexts to match the query results as they are revealed. For our version of the PRF and encryption scheme above, the simulator can arrange for the random oracle responses to point at random labels, and for the ciphertexts to decrypt to the revealed results. The only defects in the simulation occur when an adversary manages to query the random oracle with a key before it is revealed, which can be shown to happen with negligible in λ probability. \square

ALTERNATIVE APPROACH TO ADAPTIVE SECURITY. We sketch how to modify our protocol to achieve adaptive security without a random oracle at the cost of extra communication. We choose the encryption scheme for the scheme to be of the one-time pad form e.g. CTR mode with a random IV. Now instead of sending the keys K_1 and K_2 , the client computes the labels and encryption pads herself and sends them to the server, who can retrieve the labels and perform the decryption. In general the client will not know when to stop, but we can either have the client retrieve a server-stored encrypted counter first, or have the server send a "stop" message when all of the results have been found. Note that the required additional communication is proportional to the size of the result-set and can overlap the disk access as well as the return of results. Hence, the resulting scheme should perform in practice as good as the prior schemes.

ASYMPTOTIC EFFICIENCY. The EDB consists of a dictionary holding $N = \sum_{w \in W} |\text{DB}(w)|$ identifier/ciphertexts pairs. Searching is fully parallelizable if the dictionary allows parallel access, as each processor can independently compute $F(K_1, c)$ and retrieve/decrypt the corresponding ciphertext.

RELATION TO [5] A prior SSE scheme by Chase and Kamara used a dictionary, but in a crucially different way. There, a single label was associated with the entire set $\text{DB}(w)$, and thus security requires padding all of the result sets to the maximum size. We instead associate one label with each result for a keyword (so if there are T documents with a keyword, then

there are T independent labels in our scheme but only 1 label in the Chase-Kamara scheme). This allows us to avoid padding and enable parallel searching, resulting in significant storage savings and performance gains on large datasets.

A. Efficient extensions

We give a sequence of three schemes (denoted Π_{pack} , Π_{ptr} , $\Pi_{2\text{lev}}$, with names explained below) that exhibit the techniques in our most practical $\Pi_{2\text{lev}}$ construction.

REDUCING DICTIONARY RETRIEVALS: Π_{pack} . During a search for w , our basic construction performs $|\text{DB}(w)|$ retrievals from the dictionary, each with an independent and random-looking tag. Even an external-memory efficient dictionary will perform relatively poorly when the dictionary is stored on disk.

Most prior schemes suffer from this drawback. To improve locality we modify the basic construction to encrypt several identifiers in each ciphertext. Specifically, we fix a block size B , and when building the results list, we process B identifiers at a time and pack them into one ciphertext d , with the same tag. We pad the last block of identifiers up to the same length. Searching proceeds exactly as before, except the server decrypts and parses the results in blocks instead of individually. We denote this variant Π_{pack} .

This reduces the number of disk accesses from $|\text{DB}(w)|$ dictionary retrievals to $\lceil |\text{DB}(w)|/B \rceil$. We can prove security against non-adaptive or adaptive attacks under the same assumptions, but with the leakage function \mathcal{L}_B that starts by outputting $\sum_{w \in \mathcal{W}} \lceil |\text{DB}(w)|/B \rceil$ instead of $\sum_{w \in \mathcal{W}} |\text{DB}(w)|$. We note that this leakage is incomparable to the original leakage (see the discussion at the end of this section).

Theorem 4: Π_{pack} is correct and \mathcal{L}_B -secure against non-adaptive attacks if F is a secure PRF and (Enc, Dec) is RCPA-secure.

The proof is a simple extension of the proof for the basic construction. The essential observation is that the server only needs to know how many of the packed blocks to create in the encrypted index. Similar to before, we can achieve adaptive security in the random oracle model or by increasing communication. We defer the details to the full version.

FURTHER REDUCTION VIA POINTERS: Π_{ptr} . Π_{pack} would be inefficient when returning very large sets $\text{DB}(w)$, as the server still performs $\lceil |\text{DB}(w)|/B \rceil$ dictionary retrievals. Making B large results in too much padding when the dataset contains many keywords only appearing in a few $\ll B$ documents.

We address this tension by modifying the scheme again, calling the new variant Π_{ptr} . Π_{ptr} packages the identifiers into encrypted blocks of B as before, but it stores these blocks in random order in external memory and not in the dictionary (technically, we say they are stored in an array). The scheme will now use the dictionary to store encrypted blocks of b pointers to these encrypted blocks. To search, the server will retrieve the encrypted pointers from the dictionary and then follow them to the encrypted identifiers.

Π_{ptr} is described in Figure 3. In this scheme, the EDB consists of a dictionary γ holding encrypted blocks of b pointers

Setup(DB)

1. $K \xleftarrow{\$} \{0, 1\}^\lambda$; allocate array A , list L
2. For each $w \in \mathcal{W}$:
 - Set $K_i \leftarrow F(K, i||w)$ for $i = 1, 2$
 - $T \leftarrow \lceil |\text{DB}(w)|/B \rceil$
 - Partition $\text{DB}(w)$ into B -blocks I_1, \dots, I_T
 - Pad I_T up to B entries if needed
 - Choose random empty indices i_1, \dots, i_T in A
 - For $j = 0, \dots, T$ (*store id blocks in array*)
 - $d \leftarrow \text{Enc}(K_2, I_j)$; Store $A[i_j] \leftarrow d$
 - $T' \leftarrow \lceil T/b \rceil$
 - Partition $\{i_1, \dots, i_T\}$ into b -blocks $J_1, \dots, J_{T'}$
 - Pad $J_{T'}$ up to b entries if needed
 - For $c = 0, \dots, T'$ (*store ptr blocks in dict.*)
 - $\ell \leftarrow F(K_1, c)$; $d' \leftarrow \text{Enc}(K_2, J_c)$
 - Add (ℓ, d') to L
 - Set $\gamma \leftarrow \text{Create}(L)$
3. Output the client key K and EDB = (γ, A) .

Search

- Client:* On input (K, w) ,
- $K_1 \leftarrow F(K, 1||w)$, $K_2 \leftarrow F(K, 2||w)$
 - Send (K_1, K_2) to the server.
- Server:* For $c = 0$ until Get returns \perp ,
- $d \leftarrow \text{Get}(\gamma, F(K_1, c))$
 - $(i_1, \dots, i_b) \leftarrow \text{Dec}(K_2, d)$
 - For $j = 0, \dots, b$ (*ignore padding indices*)
 - $m \leftarrow \text{Dec}(K_2, A[i_j])$
 - Parse and output ids in m

Fig. 3. Scheme Π_{ptr} .

and an array A holding blocks of B encrypted identifiers for a given keyword, where b and B are block size parameters to be chosen. The setup algorithm stores blocks of encrypted results in random locations in A , and then stores encrypted pointers to those locations in γ , with labels that allow retrieval similar to the prior variants.

In the full version we show that this variant achieves the security for the leakage function $\mathcal{L}_{b,B}$ which initially outputs $\sum_{w \in \mathcal{W}} \lceil |\text{DB}(w)|/B \rceil$ and $\sum_{w \in \mathcal{W}} \lceil |\text{DB}(w)|/(bB) \rceil$, which are the number of blocks in γ and A respectively, and later leakages are just the identifiers as before.

MOST PRACTICAL VARIANT: $\Pi_{2\text{lev}}$. In real data sets the number of records matched by different keywords will vary by several orders of magnitude. This presents a challenge in optimizing our variants, and we could not find a setting of B and b that gave an acceptable trade-off between index size (due to padding) and search time. Small sets $\text{DB}(w)$ still resulted in a large block of size B in the dictionary and a large block of size b in the array, while huge sets still required many retrievals from the dictionary.

Thus we again modify the scheme to extend the ideas before, calling the new variant $\Pi_{2\text{lev}}$. The crucial difference is that sets $\text{DB}(w)$ can be processed and stored differently based on their sizes, with an extra level of indirection for very large sets that explains the name. Care must be taken to do this with an acceptable form of leakage.

Below we describe the Π_{pack} variant formally. At a high

Setup(DB)

1. $K \xleftarrow{\$} \{0, 1\}^\lambda$ allocate list L , array A
2. For each $w \in W$
 - Set $K_i \leftarrow F(K, i||w)$ for $i = 1, 2$
 - $T \leftarrow \lceil |\text{DB}(w)|/B \rceil$
 - If $|\text{DB}(w)| \leq b$ (*small case: ids in dict.*)
 - Pad $\text{DB}(w)$ to b elements
 - $\ell \leftarrow F(K_1, 0)$; $d \leftarrow \text{Enc}(K_2, \text{DB}(w))$
 - Add (ℓ, d) to L
 - If $|\text{DB}(w)| > b$ (*medium/large cases*)
 - Partition $\text{DB}(w)$ into B -blocks I_1, \dots, I_T
 - Pad I_T up to B elements
 - Choose random empty indices i_1, \dots, i_T in A
 - For $j = 1, \dots, T$ (*store ids in array*)
 - $d \leftarrow \text{Enc}(K_2, I_j)$; $A[i_j] \leftarrow d$
 - If $T \leq b$ (*medium case: ptrs in dict.*)
 - Pad $\{i_1, \dots, i_T\}$ to b elements
 - $\ell \leftarrow F(K_1, 0)$; $d' \leftarrow \text{Enc}(K_2, i_1 || \dots || i_b)$
 - Add (ℓ, d) to L
 - If $b < T \leq bB$ (*large case: ptrs in array & dict.*)
 - $T' \leftarrow \lceil T/B \rceil$
 - Part. $\{i_1, \dots, i_T\}$ into b -blocks $J_1, \dots, J_{T'}$
 - Pad $J_{T'}$ to B elements
 - Choose random empty indices $i'_1, \dots, i'_{T'}$ in A
 - For $j = 1, \dots, T'$
 - $d \leftarrow \text{Enc}(K_2, J_j)$; $A[i'_j] \leftarrow d$
 - Pad $\{i'_1, \dots, i'_{T'}\}$ to b elements
 - $\ell \leftarrow F(K_1, 0)$; $d'' \leftarrow \text{Enc}(K_2, i'_1 || \dots || i'_b)$
 - Add (ℓ, d'') to L
3. $\gamma \leftarrow \text{Create}(L)$
Output the client key K and $\text{EDB} = (\gamma, A)$.

Fig. 4. Setup for SSE Scheme $\Pi_{2\text{lev}}$.

level, it works as follows. It classifies the sets $\text{DB}(w)$ as *small*, *medium*, or *large*. For small sets, it will store the identifiers directly in the dictionary (so no pointers are used, similar to the packed variant Π_{pack}). For medium size sets, it will store them as in Π_{ptr} , with a block of pointers in the dictionary and then blocks of identifiers in the array. Finally large sets are stored differently, with two levels of indirection: The dictionary is now used to hold pointers that point to blocks of pointers in the array, which point to the identifier blocks.

In $\Pi_{2\text{lev}}$ we again fix parameters b and B to be sizes of blocks in an dictionary γ and array A respectively. The scheme classifies each of the result sets $\text{DB}(w)$ with $|\text{DB}(w)| \leq b$ as *small*, sets of size $b < |\text{DB}(w)| \leq Bb$ as *medium*, and finally sets of size $Bb \leq |\text{DB}(w)| < B^2b$ as *large*. We will always set b, B so that no set is larger than B^2b .

Small sets fit completely in a block of the top-level dictionary γ , and are stored there. Medium sets will be stored as in the previous variant but with a single block of at most b pointers in γ and the corresponding blocks of identifiers in A . These sets consist of between $b + 1$ and bB identifiers.

Finally, for large sets we store a block of at most b pointers in γ . In each of the b positions pointed to in A , we store another block of at most B pointers to other positions in A . Finally, these pointers point to blocks of encrypted identifiers. Figure 4 describes the $\text{Setup}(\text{DB})$ function of $\Pi_{2\text{lev}}$ in more detail.

To search, the client works as with the other variants by sending K_1, K_2 . The server computes the label $\ell \leftarrow F(K_1, 0)$, and retrieves $d \leftarrow \text{Get}(\gamma, \ell)$, and decrypts d using K_2 . If it

finds identifiers here, then it outputs them and stops. Otherwise, it uses the pointers to retrieve blocks from A . If those blocks contain identifiers then it outputs them. Otherwise it follows the next level of pointers to finally find the identifiers, which it decrypts and outputs.

SECURITY. We prove security for the leakage function $\mathcal{L}_{m,b,B}$ that initially outputs $m = |W|$ and the value

$$S = \sum_{w:|\text{DB}(w)|>b} \lceil |\text{DB}(w)|/B \rceil + \sum_{w:|\text{DB}(w)|>bB} \lceil |\text{DB}(w)|/B^2 \rceil.$$

The value m is the number of data items in γ , and the value S is the number of blocks in A . This is leaking S itself, which is defined by the above sum, and not the individual summands, resulting leakage that is incomparable to our other variants and to prior schemes. On search queries \mathcal{L} has the same leakage as before.

Theorem 5: $\Pi_{2\text{lev}}$ is correct and $\mathcal{L}_{m,b,B}$ -secure against non-adaptive attacks if F is a secure PRF and (Enc, Dec) is RCPA-secure.

We can prove adaptive security either in the random oracle model or by increasing communication. We defer this simple extension to the full version.

POINTERS VS. IDENTIFIERS. Although pointers are smaller than identifiers in our implementations, $\Pi_{2\text{lev}}$ packs the same number of pointers or identifiers together (b in the dictionary, or B in the array) to simplify the presentation. The actual implementation packs more pointers into a block than identifiers. Formally, we introduce parameters b', B' , and modify $\Pi_{2\text{lev}}$ as follows.

- When storing identifiers in the dictionary (in the small case), it packs up to b of them together, but when storing pointers there it packs b' in the same amount of space.
- When storing identifiers in the array (in the medium and large cases), it packs up to B of them together, but when storing pointers there it packs B' together in the same amount of space.

This causes an analogous change to the value S leaked, which can be calculated similarly. We defer the formal analysis (which is almost identical to that of $\Pi_{2\text{lev}}$) to the full version.

LEAKAGE DISCUSSION. The leakage functions $\mathcal{L}_B, \mathcal{L}_{b,B}, \mathcal{L}_{m,b,B}$ are non-standard. First consider \mathcal{L}_B , and how it compares to \mathcal{L} which outputs $N = \sum_{w \in W} |\text{DB}(w)|$. Any input DB with m unique keywords, each with $|\text{DB}(w)| \leq b$, will be indistinguishable under \mathcal{L}_B , but many of them will not be under \mathcal{L} . A similar incomparability goes in the other direction. We are not aware of a scenario where this difference is important for reasonably small B . The function $\mathcal{L}_{b,B}$ leaks strictly more information than \mathcal{L}_B (actually \mathcal{L}_b), but it also does not appear to be harmful. Finally, $\mathcal{L}_{m,b,B}$ leaks this type of size information *and* the number of keywords m . The number m seems to be the most useful information for an adversary, but in prior work it has been considered acceptable. It is possible to modify the scheme to avoid leaking exactly m , say by storing blocks of a different size in the dictionary.

IV. DYNAMIC CONSTRUCTIONS

We extend our static SSE constructions to support changes to the database. Our dynamic SSE constructions will consist of a statically encrypted database EDB using any of the schemes described above, and an auxiliary encrypted database EDB⁺ which is maintained to be of the form of a basic dictionary-based scheme. The EDB⁺ is initially empty and changes as updates happen.

ADD-ONLY SCHEME: Π_{bas}^+ . We start with an extension of Π_{bas} , denoted Π_{bas}^+ that supports additions only, meaning $\text{add}, \text{edit}^+$ inputs from the client during Update. Π_{bas}^+ is simpler and possibly interesting in its own right.

To support additions we use a dictionary γ^+ which is initially empty and to which a pair (ℓ, d) is added with each keyword addition; here ℓ is a label computed from the keyword and a keyword-specific counter, and d is the encryption of the record id involved in the addition operation. Search for a keyword w is performed by the server by first searching γ as in the static case, then re-computing all labels corresponding to w in γ^+ . The latter labels are computed using a w -specific key provided by the client and a running counter.

Note that addition operations involving keyword w require the client to know the current value of the w -specific counter. For this, the scheme maintains a dictionary δ associating each keyword that was ever added via edit^+ or add with its current counter value. δ can be stored at the client or stored at the server and retrieved by the client for performing update operations. We formalize a scheme Π_{bas}^+ where the client stores locally the dictionary δ and discuss below a stateless variant. We assume throughout that the client never tries to add a record/keyword pair that is already present - it is easy, but messy, to extend our scheme and the leakage profile to handle this.

In Π_{bas}^+ , $\text{Setup}(\text{DB})$ is exactly as in Π_{bas} except that the client also initializes δ to be an empty dictionary and keeps it as state, and the server initializes an empty dictionary γ^+ that is stored with EDB. We next give the update protocol.

Update: We only specify the protocol with client input $\text{op} \in \{\text{add}, \text{edit}^+\}$. The parties work exactly the same on either type of operation. The client has input id, W_{id} . It starts by deriving key $K^+ \leftarrow F(K, 3)^2$, and proceeds as follows:

For $w \in W_{\text{id}}$:

- Set $K_1^+ \leftarrow F(K^+, 1||w)$, $K_2^+ \leftarrow F(K^+, 2||w)$.
- $c \leftarrow \text{Get}(\delta, w)$; If $c = \perp$ then $c \leftarrow 0$
- Set $\ell \leftarrow F(K_1^+, c)$; $d \leftarrow \text{Enc}(K_2^+, \text{id})$
- $c++$; Insert (w, c) into δ
- Add (ℓ, d) to L in lexicographic order
- Send L to the server.

When inserting (w, c) into δ , we assume that it will overwrite any previous entry (w, \cdot) if it exists.

Finally, the server adds each $(\ell, d) \in L$ to γ^+ . This completes the update protocol.

To complete Π_{bas}^+ we describe the protocol Search.

Search: On input w , the client sets $K^+ \leftarrow F(K, 3)$ and proceeds:

- $K_1 \leftarrow F(K, 1||w)$, $K_2 \leftarrow F(K, 2||w)$
- $K_1^+ \leftarrow F(K^+, 1||w)$, $K_2^+ \leftarrow F(K^+, 2||w)$
- Send (K_1, K_2, K_1^+, K_2^+) to the server.

Upon receiving the message, the server computes its output as follows:

- For $c = 0$ until Get returns \perp ,
 - $d \leftarrow \text{Get}(\gamma, F(K_1, c))$; $m \leftarrow \text{Dec}(K_2, d)$
 - Parse and output id in each m
- For $c = 0$ until Get returns \perp ,
 - $d \leftarrow \text{Get}(\gamma^+, F(K_1^+, c))$; $m \leftarrow \text{Dec}(K_2^+, d)$
 - Parse and output id in each m .

Intuitively, the server is repeating the search procedure from Π_{bas} twice: Once with (K_1, K_2) and γ , and then with (K_1^+, K_2^+) and γ^+ .

LEAKAGE PROFILE FOR Π_{bas}^+ . Let us first give some intuition for the leakage of Π_{bas}^+ . Initially the leakage is exactly like Π_{bas} , where only the size of DB is leaked. Upon an edit^+ or add query, if the keywords being added were not previously searched for, then the server learns nothing other than number of record/keyword pairs added (not even the if the operation was edit^+ vs. add). If, however, one (or more) of the keywords were previously searched for, then the server can reuse its keys from before to detect the presence of these keywords (this type of leakage is inherent when the keys provided to the server for searching are deterministically generated and the same each time). The leakage on a search is similar to before, except now for record/keyword pairs in γ^+ the server can recognize when they were added. The order for pairs in γ generated at setup time is still hidden, however.

We proceed with the formal definition of \mathcal{L}^+ for adaptive security. Amongst its state, it will keep a list Q describing all queries issued so far, where an entry of Q is of the form (i, op, \dots) , meaning a counter, the operation type, and then the one or more inputs to the operation.

On initial input DB, \mathcal{L}^+ creates a state consisting of a counter $i \leftarrow 0$, an empty list Q and DB, and a set ID initialized to contain all of the identifiers in DB. Let us define the *search pattern* $\text{sp}(w, Q)$ of a keyword with respect to Q to be the indices of queries that searched for the keyword w , i.e.

$$\text{sp}(w, Q) = \{j : (j, \text{srch}, w) \in Q\}.$$

For an identifier id and keyword w , the *add pattern* of id, w with respect to Q is the indices that added w to the document id , i.e.

$$\begin{aligned} \text{ap}(\text{id}, w, Q) = & \{j : (j, \text{add}, \text{id}, W_{\text{id}}) \in Q, w \in W_{\text{id}}\} \\ & \cup \{j : (j, \text{edit}^+, \text{id}, W_{\text{id}}) \in Q, w \in W_{\text{id}}\}. \end{aligned}$$

Finally, we let the *add pattern of keyword w with respect to Q and ID* be the set of all identifiers to which w was ever added (via a add or edit^+ operation) along with the indices showing when they were added. That is,

$$\text{AP}(w, Q, \text{ID}) = \{(\text{id}, \text{ap}(\text{id}, w, Q)) : \text{id} \in \text{ID}, \text{ap}(\text{id}, w, Q) \neq \emptyset\}.$$

²We use input 3 for domain separation to K^+ make independent of each keyword-specific $K_1 = F(K, 1||w)$ and $K_2 = F(K, 2||w)$.

\mathcal{L}^+ produces outputs for the initial query, edit^+ and add updates, and search queries as follows:

- On initial input DB it saves state as defined above and outputs $N = \sum_{w \in W} |\text{DB}(w)|$.
- For a search query w , \mathcal{L}^+ appends (i, srch, w) to Q and increments i . Then it outputs $\text{sp}(w, Q)$, $\text{DB}(w)$, and $\text{AP}(w, Q, \text{ID})$.
- Update queries for edit^+ and add operations are handled similarly. For a query $(\text{edit}^+/\text{add}, \text{id}, W_{\text{id}})$, \mathcal{L}^+ first appends $(i, \text{edit}^+/\text{add}, \text{id}, W_{\text{id}})$ to Q , adds id to ID , and increments i . It outputs $|W_{\text{id}}|$ and the (lexicographically ordered) set of search patterns

$$\{\text{sp}(w, Q) : w \in W_{\text{id}}\}.$$

If any of the search patterns was non-empty, then it also outputs id .

While subtle in its formulation, \mathcal{L}^+ is essentially the best possible leakage for an SSE scheme that generates the same search keys on repeated searches.

In words, the search query leakage includes $\text{sp}(w, Q)$ and $\text{DB}(w)$ for obvious reasons. The add pattern of w , $\text{AP}(w, Q, \text{ID})$, is the set of id matching w added later along with “history” information $\text{ap}(\text{id}, w, Q)$ indicating when they added. The order information represents that the server can look at γ^+ and see when each id was added by rewinding and re-running searches. For updates Π_{bas}^+ leaks *only the size of the update* if the added keywords have not been searched for. If any of them have been searched for, then the server learns that “a keyword with search pattern $\text{sp}(w, Q)$ was added” via the set of search patterns in the update leakage. Finally it learns the id being updated because it has the ability to search for any of its keywords. Each of these leakage components is unavoidable for a deterministic SSE scheme, and we regard them as minimal.

We can now state our security for the add-only scheme. A proof will appear in the full version.

Theorem 6: Π_{bas}^+ is correct and \mathcal{L}^+ -secure against non-adaptive attacks if F is a secure PRF and (Enc, Dec) is RCPA-secure.

STATELESS CLIENT VARIANT. Above, the client keeps a dictionary δ containing one counter per keyword that is added after initialization. We could modify the scheme so that the client is stateless by storing δ in encrypted form at the server and having the client download and re-encrypt *all of* δ for each update (note that the size of δ is as the number of distinct keywords added via add and edit^+ and *not* the total number of keywords in the set W). In this variant the server will learn how many *new* keywords are added each time by watching if δ grows. We leave a formal version of this scheme to the full version.

DYNAMIC SCHEME $\Pi_{\text{bas}}^{\text{dyn}}$. We now augment the Π_{bas} scheme with del and edit^- operations to obtain our fully dynamic scheme $\Pi_{\text{bas}}^{\text{dyn}}$. We will implement deletions by maintaining a revocation list and having the server discard results that have been revoked.

To delete a record/keyword pair (id, w) from the server’s storage, the client will generate a pseudorandom *revocation identifier* and send it to the server. During searches, the client will give the server a key that allows it to recompute revocation identifiers, which it will then use to filter out deleted results. This complicates our addition protocols. To add a pair that was previously deleted, the protocol must “unrevoke” that pair by having the server delete its revocation identifier.

We now formally specify $\Pi_{\text{bas}}^{\text{dyn}}$. Setup is exactly the same as Π_{bas}^+ , except that the server also initializes an empty set S_{rev} . As a data structure, S_{rev} will support additions, deletions, and membership testing.

Update: We first describe how to handle client inputs with $\text{op} \in \{\text{del}, \text{edit}^-\}$. The client takes as input $(\text{del}/\text{edit}^-, \text{id}, W_{\text{id}})$, and first derives a key $K^- = F(K, 4)$ ³, and then computes

For $w \in W_{\text{id}}$ do

$$K_1^- \leftarrow F(K^-, w), \text{revid} \leftarrow F(K_1^-, \text{id})$$

Add revid to L_{rev} in lexicographic order

Send L_{rev} to the server.

The server receives L_{rev} and adds each revid to S_{rev} . This completes Update for the del and edit^- operations.

Next we define Update for $\text{op} \in \{\text{add}, \text{edit}^+\}$. On input $(\text{add}/\text{edit}^+, \text{id}, W_{\text{id}})$, the client performs a computation similar to the list L computation in Π_{bas}^+ , except that it also includes the appropriate revid values. It then awaits a response from the server specifying which additions resulted in a true addition and which caused an “unrevocation”, and uses this information to increment the correct counters.

In code, the client sets $K^- \leftarrow F(K, 4)$ and does the following:

For $w \in W_{\text{id}}$:

$$K_1^+ \leftarrow F(K^+, 1||w); K_2^+ \leftarrow F(K^+, 2||w)$$

$$K_1^- \leftarrow F(K^-, w)$$

$$c \leftarrow \text{Get}(\delta, w); \text{If } c = \perp \text{ then } c \leftarrow 0$$

$$\ell \leftarrow F(K_1^+, c); d \leftarrow \text{Enc}(K_2^+, \text{id})$$

$$\text{revid} \leftarrow F(K_1^-, \text{id})$$

Add (ℓ, d, revid) to L in lexicographic order

Send L to the server.

The server generates its response $r \in \{0, 1\}^{|L|}$ as follows. For the i -th pair $(\ell, d, \text{revid}) \in L$ in order, if $\text{revid} \in S_{\text{rev}}$, it sets the i -th bit of r to 1 and deletes revid from S_{rev} . Else, it clears that bit to 0 and adds (ℓ, d) to γ . Finally, it sends r to the client.

Now the client increments the counters for keywords corresponding to 0 bits in r . It processes the keywords $w \in W_{\text{id}}$ in order of their labels in L . For the i -th keyword w in that order, if the i -th bit of r is 0 it computes $c \leftarrow \text{Get}(\delta, w)$, increments c , and inserts (w, c) into δ . This completes the update protocol.

The last component of $\Pi_{\text{bas}}^{\text{dyn}}$ is the search protocol.

Search: On client input w , it sets $K^- = F(K, 4)$, $K_1^- = F(K^-, w)$, and then computes (K_1, K_2, K_1^+, K_2^+) as in Π_{bas}^+ .

³As with K^+ the input 4 is for domain separation only.

It sends $(K_1, K_2, K_1^+, K_2^+, K_1^-)$ to the server. The server computes the result identifiers using the first four keys exactly as in Π_{bas}^+ , except before outputting each id it computes $\text{revid} = F(K_1^-, \text{id})$ and tests if $\text{revid} \in S_{\text{rev}}$. If so, it discards id instead of outputting it.

LEAKAGE FUNCTION. We now define the leakage profile \mathcal{L}_{dyn} . It will maintain a list of query information Q and set of identifiers ID like \mathcal{L}^+ from above. Below we use the same definitions for $\text{sp}, \text{ap}, \text{AP}$ as in \mathcal{L}^+ , and define the following analogous patterns dp, DP for deletions:

$$\begin{aligned} \text{dp}(\text{id}, w, Q) &= \{j : (j, \text{del}, \text{id}, W_{\text{id}}) \in Q, w \in W_{\text{id}}\} \\ &\cup \{j : (j, \text{edit}^-, \text{id}, W_{\text{id}}) \in Q, w \in W_{\text{id}}\}. \end{aligned}$$

and $\text{DP}(w, Q, \text{ID}) =$

$$\{(\text{id}, \text{dp}(\text{id}, w, Q)) : \text{id} \in \text{ID}, \text{dp}(\text{id}, w, Q) \neq \emptyset\}.$$

Intuitively, $\text{dp}(\text{id}, w, Q)$ is the set of indices of queries that deleted w from id , and $\text{DP}(w, Q, \text{ID})$ is the set of identifiers from which w was deleted, along with the corresponding deletion pattern.

- On first input DB , \mathcal{L}_{dyn} initializes a counter $i \leftarrow 0$, empty list Q , set ID to be identifiers in DB . It saves $\text{DB}, i, \text{ID}, Q$ as state, and outputs $N = \sum_{w \in \mathcal{W}} |\text{DB}(w)|$.
- On search input w , \mathcal{L}_{dyn} appends (i, srch, w) to Q , increments i , and outputs $\text{sp}(w, Q)$, $\text{DB}(w)$, $\text{AP}(w, Q, \text{ID})$, and $\text{DP}(w, Q, \text{ID})$.
- On update input $(\text{add}/\text{edit}^+, \text{id}, W_{\text{id}})$, it appends $(i, \text{add}/\text{edit}^+, \text{id}, W_{\text{id}})$ to Q , adds id to ID , and increments i . It outputs $\text{add}, |W_{\text{id}}|$ and the set $\{(\text{sp}(w, Q), \text{ap}(\text{id}, w, Q), \text{dp}(\text{id}, w, Q)) : w \in W_{\text{id}}\}$. Finally, if any of the $\text{sp}(w, Q)$ are non-empty, then it also outputs id .
- On update input $(\text{del}/\text{edit}^-, \text{id}, W_{\text{id}})$, it appends $(i, \text{del}/\text{edit}^-, \text{id}, W_{\text{id}})$ to Q , adds id to ID , and increments i . Then it computes its output exactly as in the add/edit^+ case above, except that it outputs del instead of add as the first component.

The leakage on searches is minimal: It consists of all patterns of searches, deletions, and additions that can be derived once the server has the ability to search for a keyword and rewind the database. For leakage on updates, the server will learn when/if that identifier has had the same keywords added or deleted before, and also when/if the same keywords have been searched for. This comes from observing the revid values, which will repeat every time the same identifier/keyword pair is added or deleted. Note that, if same keyword is added/deleted from two documents, then this information is not leaked until it is searched for (contrast this with [14] which leaks this information always).

We prove the following theorem, as well as an adaptive variant, in the full version.

Theorem 7: $\Pi_{\text{bas}}^{\text{dyn}}$ is correct and \mathcal{L}_{dyn} -secure against non-adaptive attacks if F is a secure PRF and (Enc, Dec) is RCPA-secure.

ASYMPTOTIC ANALYSIS. To add a file the client sends one label/ciphertext/revid per record/keyword pair being changed. For deletions, the δ dictionary is not involved. The client just sends one revid per document/keyword to be deleted. Assuming the dictionaries γ, γ^+ , and the revocation list are fully read-parallel, and the number of deletions is much smaller than the size of the EDB, each search operation continues to have the same order run-time complexity as in the basic static construction of Figure 2.

DISCUSSION AND COMPARISON TO PRIOR WORK. Our scheme $\Pi_{\text{bas}}^{\text{dyn}}$ is unsatisfying in some situations as it does not reclaim space after deletions. While this is a drawback, all known dynamic SSE schemes [14], [15], [21] have severe drawbacks in different dimensions, and no scheme has achieved an ideal combination of leakage, index size, and full functionality like reclaiming space.

The scheme of [21] has no security proof, and the scheme of [14] has a worst-case quadratic size encrypted index. The dynamic scheme in [15] has much more leakage than our scheme, effectively leaking the pattern of all intersections of everything that is added or deleted, whether or not the keywords were searched for. For an example, suppose $\{w_1, w_2, w_3\}$ are added to id_1 , $\{w_1, w_2\}$ are added to id_2 , and $\{w_1\}$ is added to id_3 . Then [15] will leak that exactly one common keyword was added to all three and that exactly two common keywords were added to the first two (but not the third) and so on. This structural “equality pattern” is the sort of leakage that we do not leak.

Not reclaiming space allows our implementations to be much simpler and also gives us the flexibility to apply various efficiency optimizations (as in section III A) to the static scheme which seem hard to achieve when in-place updates have to be supported. As our data structures are more compact than prior work, the overall space requirements will be lower anyway for some number of deletes. In particular, as compared to prior work [14] we are not forced to estimate an upper bound (by necessity, conservative) on the maximum database size.

In some settings where SSE is used as a component, the encrypted database is re-encrypted for security reasons [13]. In these settings we can reclaim space and combine the auxiliary data structure with the main static data structure while re-encrypting.

APPLICATION TO $\Pi_{\text{ptr}}, \Pi_{\text{pack}}, \Pi_{\text{2lev}}$. The dynamic extensions to Π_{bas} can be applied as-is to other variants, resulting in almost the same leakage \mathcal{L}_{dyn} . The only difference is the size leakage in the initial input DB , which changes according to the different schemes. In our implementation in the next section we consider these variants.

V. IMPLEMENTATION

We report on our implementations of Π_{2lev} and Π_{pack} (described in Section III), with extensions for dynamic data updates (Section IV). The former scheme is the most efficient and scales to the largest datasets; it represents our current prototype. The latter is a simplification of the original OXT implementation which we introduced in [3] and is discussed here to better illustrate the effectiveness of the ideas in Π_{2lev} and the improvement over prior work.

PRACTICAL CRITERIA. Before describing our results, we enumerate some of the practical criteria that we optimize for in the Π_{2lev} prototype.

- **Parallel EDB access:** The server should be able to issue concurrent access requests to EDB when processing a search. Modern disk controllers handle thousands of concurrent requests and optimize disk access patterns, increasing transfer bandwidth by orders of magnitude when compared with sequential access. Requests are served out-of-order but the performance benefits offset the additional implementation complexity.
- **EDB goodput:** EDB design should maximize I/O goodput, i.e., the ratio of data used by the processing of a query relative to the total amount of data retrieved from external storage. In addition to selecting an appropriate dictionary, we achieve this by setting the parameters b, b', B, B' in Π_{2lev} to take maximum advantage of the block device.
- **Small EDB storage:** The dictionary used in EDB should minimize storage overhead while satisfying the other constraints.
- **Lightweight EDB updates:** Update information will be independent from the EDB and implemented in-memory. This is consistent with our envisioned scenarios where updates are either infrequent or periodically folded into the main data structure via re-encryption of the entire database.

INPUT DATASETS. Our implementation accepts as input both relational databases and document collections. The latter are mapped to relational database tables with document attributes, such as author name, creation date, etc., stored in atomic columns and with the document content stored in a text column.

We target clear-text datasets (DBs) that consist of several tens of billions of distinct (keyword, id) pairs. The EDBs generated from such datasets take several terabytes of storage and require several times more temp storage for Setup. We aim to process such datasets efficiently (Setup(DB) and Search) on medium size 64-bit x86 platforms (in our configuration, 8 cores, 96GB of RAM, and ≈ 100 TB RAID volume on external storage box).

The constructions described in this paper and their implementations can be extended to support richer functional settings than simple keyword search, such as SSE in multi-client settings or boolean queries via the OXT protocol [3] (see end of Section II), by storing in the EDB for each (keyword, id) pair more data than just the encrypted document id. In the following, we use the term *tuple* for the data stored per (keyword, id) pair in any of these functional settings.

ORGANIZATION. The next two subsections describe our experiences with the Π_{pack} prototype, which is the subset of the OXT implementation [3] relevant to this work, and the design and implementation of our Π_{2lev} (see Figure 4). A particular challenging issue for both prototypes was EDB generation time; the Setup implementation for Π_{2lev} is discussed separately in Section V-C. Section V-D describes how these

constructs are used to support richer functional settings, such as OXT. Finally, Section V-E describes several representative experiments.

A. Π_{pack} Implementation

The discussion of the Π_{pack} implementation here is intended as a preamble to our presentation of Π_{2lev} in the next subsection as it serves to motivate the optimizations applied to the latter construction. Our implementation of Π_{pack} instantiates the EDB dictionary using a bucket hash table. Buckets are split in equal-size *locations*, which are used to store equal-size groups of tuples created by partitioning the $DB(w)$ sets. The location size is equal to the group size plus the size of the attached label. Each group can be stored in any of the free locations in the bucket determined by hashing its label. As usual, the hash map is over-sized to allow for all groups to be placed successfully; empty locations are filled with random bits to mask the total number of groups in the EDB.

Using a bucket hash for the dictionary allowed us to avoid sorting the tuples by label (as required for security) before creating the dictionary. This worked by ensuring the dictionary is *history independent*, meaning the output of $Create(L)$ depends only on the members of L and not on the order they were added to L .

The bucket hash table is stored in one large file on an ext4 RAID partition of attached storage. The bucket size is set to a multiple of the RAID stripe size⁴, and buckets are aligned with the pages of the underlying file system.

The two most significant drawbacks with the above construction are the need for oversizing the hash table, which translates into a much larger EDB than needed, and the poor goodput, as one has to retrieve an entire bucket to access a group of tuples. In experiments with configurations and data sets similar to those described in [3], the hash table has a load factor of $\approx 60\%$ (i.e., over-sized by a factor of ≈ 1.6) for the placement to be successful, and goodput is $\approx 1\%$, as there are 96 locations per bucket.

To achieve a higher load factor (smaller EDB), we built another Π_{pack} prototype which uses a Cuckoo Hash (CH) table modeled after [8] for the dictionary; page size and alignment are the same as for the bucket hash dictionary in the previous construction. Although we achieve load factors a little over 90%, the cost of handling collisions during EDB generation is very high. Moreover, making the dictionary *history independent* is much more difficult when using a CH table and likely inefficient in our setting.

We designed a more efficient algorithm to handle collisions during EDB generation, which leverages the server memory, but we found its performance to be limited by its database access patterns (see Section V-E). Finally, the need to improve the goodput motivated the design of Π_{ptr} and Π_{2lev} .

B. Π_{2lev} Implementation

In order to meet the criteria stated at the beginning of this section and avoid the drawbacks of Π_{pack} , we developed the

⁴Stripe is the smallest amount of data that can be addressed within the RAID. This is functionally equivalent to a block for an individual disk.

Π_{2lev} construction (see Figure 4) which uses different database patterns to speed-up Setup, can be configured to run Setup efficiently on platforms with a wide range of internal memory, and supports much faster retrieval as a result of higher goodput.

Recall that in Π_{2lev} , the EDB consists of two data structures: a dictionary γ and an array A . The dictionary is again implemented as a bucket hash, but now with exactly one labeled location per keyword w . The bucket address and location label are derived from w , but the location within the bucket is selected at random to ensure history independence. A γ location stores up to b tuples or b' pointers, i.e. indices in array A .

The second data structure is the array A whose entries are called *tuple blocks*. Setup uses *tuple blocks* to store tuples, or tuples and pointers, for medium or large $DB(w)$, respectively. Each *tuple block* stores up to B tuples or B' pointers, with $B \gg b$ and $B' \gg b'$ in most settings. In contrast to the dictionary γ , which is a bucket hash, the array A needs not be over-sized except for the purpose of masking the total number of tuples in EDB. Empty locations in γ and A , if any, are filled with random bits.

For all w with more than $|DB(w)| > b$, the *tuple blocks* used for $DB(w)$ are allocated at random in the array using an AES-based pseudorandom permutation and the tuple list in $DB(w)$ is split into *tuple blocks* (see medium/large cases in Figure 4). For any given w , if the number of *tuple blocks* needed to store $DB(w)$ is larger than the number of pointers that fit in a dictionary location, we use additional *tuple blocks* to store pointers (see large case Figure 4).

The dictionary γ and the array A are realized as two separate files on the same or separate ext4 RAID partitions. The location, bucket and *tuple block* sizes are configurable, but for efficiency the bucket and *tuple block* sizes must be a multiple of the RAID stripe size. Similarly, buckets and *tuple blocks* must be aligned with the pages of the underlying file system.

In our experiments, we use a low single digit number of tuples per location and 32KB or 64KB for buckets and *tuple blocks*. Pointers are 3 or 4 bytes long, depending on the size of the array A , and tuples are between 16 and 91 bytes long, depending on the functional setting. For the document collections and relational databases in our experiments, the dictionary is between one and two orders of magnitude smaller than the array.

Unpadded, the size of the dictionary leaks the approximate number of keywords while the size of the array leaks the approximate number of EDB tuples. Therefore, masking the dictionary size, which is sensitive in many scenarios, is inexpensive given its relative small size. Leaking the total number of tuples is less sensitive, which means that the larger data structure requires less or no padding in most common cases.

This construction has several important advantages for very large datasets, in particular for those having multi-modal distributions, e.g., some $DB(w)$ sets that are very large and a very large number of very small $DB(w)$ as commonly encountered. For instance, for datasets of tens of millions of documents, each English word that is not rare can be found in millions or more documents. On the other hand, ISBN

or SSN values are associated with only one book or person, respectively, independent of how many books or persons the dataset describes.

Π_{2lev} can be configured to be disk-efficient in both extremes. For rare keywords, configurations with small location sizes, corresponding to a low single digit number of tuples, allow the search engine to retrieve all the relevant tuples with only *one* disk access. Using a small location size helps reduce the dictionary size, potentially to less than the amount of the available RAM.

For the rest of the keywords, after *one* access (or a few disk accesses for very common keywords), the addresses of all the relevant *tuple blocks* are known. At this point, the query execution engine issues as many concurrent *tuple block* requests as the RAID controller can handle. After less than the average disk access time, because of the large number of pending requests, *tuple blocks* are read at close to the maximum rate of the RAID controller. The rate at which tuples are retrieved from storage determines the throughput of search engine. Note that goodput is 100% when accessing *tuple blocks* filled with tuples and that for frequent keywords, the goodput of a Search operation grows asymptotically to 100%.

In contrast, and by way of comparison, the Π_{pack} construction computes the location addresses of all their tuple groups from the keyword value and a running counter. Thus it can precompute a large number of group addresses and issue requests for tuple groups immediately, i.e. no additional disk accesses to retrieve pointers are needed. But without a priori knowledge of $DB(w)$ size, which is the common case, Π_{pack} issues many more requests than necessary. Even worse, to achieve the lowest access latency for a CH-based construction, one always needs to issue two requests per expected tuple group, as the group can be stored in two positions (pages) in the CH table. Finally, these disk accesses have low goodput as each bucket contains multiple tuple groups. Thus it appears that low I/O goodput is inherent to Π_{pack} . For large sets, the superior goodput of our construction (due to large *tuple blocks*) more than compensates for the extra initial storage access(es) for pointers.

For keywords with just a few tuples that fit in one dictionary location, the performance is the same. However, one could accelerate the performance of Π_{2lev} by storing the dictionary, which is relatively small even for large data sets, in main memory. Dictionaries used by previous work, which use one large bucket hash for all tuple sets, are too large for this optimization.

The two-level Π_{2lev} construction allows for a *very efficient* EDB generation process. As an example, during the longest phase of EDB generation from a database with ≈ 25 billion (w, id) pairs in the context of multi client OXT [3], which took 40 hours, all cores performed crypto operations at close to 100% utilization while at the same time reading 100 million records from a MySQL DB and writing to the file system the *tuple blocks* and the temp dictionary files. Overall, the two-level construction is much closer to our requirements than any previous ones and the experimental results confirm our expectations.

C. EDB Generation

For our largest datasets, EDB is on the order of 2TB. Thus EDB generation time is sensitive to implementation details and is the dominant factor determining the practical scalability of all our constructions. This section describes the parallel Setup algorithm used in the Π_{2lev} prototype.

Before EDB generation starts we process the input files into an index of the form expected by Π_{2lev} . For each text column 't' in the clear-text database table create a new 'text_t' table with columns *ind* and *word*. For each clear-text record and for each word 'xxxx' in its text column, add the pair (id, xxxx) to 'text_t', where id is the unique identifier of the clear-text record. The number of pairs in 'text_t' is equal to the number of clear-text records multiplied by the average number of words in column 't'. At the end, we create an index on the column 'word', which is used during Setup to retrieve $DB(w)$ for all $w = (t, xxxx)$, where 'xxxx' is a word in column 't'.

Unfortunately, for our largest databases, 'table_t' is too large for the index to fit in RAM, which makes building the index impractical. To overcome this obstacle, for each text column 't' we create multiple tables 'text_t_nn', such that (1) id-word pairs are somewhat evenly distributed across the new tables, (2) all the pairs with the same 'word' value are in the same table, and (3) the new tables are small enough for their indexes to be built efficiently on our platforms. Note that the atomic columns of the original table can undergo a similar transformation if the number of records in the clear-text table is too large for indexing.

EDB is generated in three phases. The first phase counts the number of distinct keywords w_i in the clear-text table and other statistics needed for sizing the dictionary γ and array A (or for masking the sizes of these data structures if so desired). This phase uses full-index scans and takes a small fraction of the total EDB generation time.

For performance reasons, the dictionary γ , realized as a bucket hash, is partitioned in equal size groups of consecutive buckets and its generation is split across the second and third phases. The *tuple block* array A is fully generated in the next phase.

The second phase generates the tuples in $DB(w)$, for all keywords $w = (i, val)$ using full-index scans on atomic column i . For each text column 't', the newly created 'text_t_nn' tables are scanned. Columns are processed in parallel worker threads, with approximately 1.5 workers per CPU core to hide database access latencies. For each value val such that $w = (i, val)$, the thread processing column i retrieves the all the ids corresponding to the records with val in column i and applies a random permutation to the resultant id sequence (i.e., $DB(w)$). For each id in the permuted sequence, the worker generates tuple elements with the encrypted id (and the additional tuple values rdk and y when implementing the more advanced features of the OXT protocol from [3]).

During this phase, each worker thread creates one temp file per dictionary partition in which it stores the content of the locations (tuples or pointers) assigned to any of the buckets in the partition. For better performance, the location content is buffered and appended to the partition file in 64KB data chunks. At the same time, for medium and large $DB(w)$, the

worker threads create all the necessary *tuple blocks* in the array A (see Figure 4).

During the third phase, the dictionary γ is created from the partition files generated in the previous phase. Each partition is constructed by a separate worker thread. Each worker thread reads the files generated by phase-two workers for its partition, merges their contents and creates the label and content of each dictionary entry in the partition. Next, for each bucket in its partition, the worker assigns the dictionary entries to random locations in the bucket, fills in empty locations with random bits, and writes the bucket to disk. For a typical census table, the dictionary file is almost two orders of magnitude smaller than the tuple block file.

Note that for the creation of the dictionary, the file system is accessed using append calls (to temp files in the second phase) and sequential read calls (from the temp files in the third phase), which are the most efficient ways to access large files.

However, worker threads still issue a large number of random write calls during the second phase, with one call for each *tuple block* generated. To reduce the disk head movements associated with these writes, we extended the parallel EDB generation algorithm to confine concurrent workers to a *tuple block window* sliding across the array. As a result, *tuple blocks* that are closely located on the disk are generated near simultaneously. This optimization reduces seek overheads noticeably for our largest EDBs.

During the third phase, threads issue another set of random writes when writing the dictionary buckets. These disk movements generated by these writes do not represent a major bottleneck because these writes are already clustered to the bucket hash partitions under constructions, which we always select in increasing order.

D. Complex Functional Settings

As already mentioned at the end of Section II, our constructions can be used to support richer encrypted-search settings than SSE, such as those in [3], [13]. In particular, all (single-keyword) SSE schemes presented here can be used, almost 'out-of-the-box', to instantiate the "TSet functionality" underlying the OXT protocol in the above works. The main change is on the size of tuples that increases in order to accommodate additional OXT information such as the $xind$ and y values (see Section 3.2 of [3]), and the key to decrypt the actual documents (as required in multi-client settings [13]).

Storing larger tuples requires minor configuration changes but no alteration of the original construct. More specifically, hash table buckets need to be made large enough to accommodate enough entries for all the tuples to be inserted into the table with high enough probability, i.e., without any of the buckets overflowing.

Another challenge in more complex protocols, such as OXT, is for the server to efficiently perform a two party computation which takes in-order generated data by the client and out-of-order the tuples, as retrieved from the disk by the Π_{pack} or Π_{2lev} prototypes. Maximizing the throughput of such a computation requires using complex buffer management

DB Name	Records	(w, id) pairs	EDB size
CW-MC-OXT-1	408,450	143,488,496	69.6 GB
CW-MC-OXT-2	1,001,695	316,560,959	99.8 GB
CW-MC-OXT-3	3,362,993	1,084,855,372	242.4 GB
CW-MC-OXT-4	13,284,801	2,732,311,945	903.9 GB
LL-MC-SKS-1	100,000	114,482,724	15.0 GB
LL-MC-SKS-2	1,000,000	1,145,547,173	52.0 GB
LL-MC-SKS-3	10,000,000	11,465,515,716	394.0 GB
LL-MC-SKS-4	100,000,000	114,633,641,708	3,961.3 GB

TABLE I. DATABASES

algorithms that optimize the use of available RAM between tokens and *tuple block* buffers.

E. Experimental Results

Our prototype implementation measures roughly 65k lines of C code, including test programs. Measurements reported here were performed on blades with dual Intel Xeon 2.4GHz E5620 processors having 4 cores each and running Linux. Storage consists of 6 1TB SAS disks configured as a RAID-0 with a 64KB stripe and attached via a 3 Gb/s to an LSI 1064e SAN controller and formatted as an ext4 file system with an 8KB page size. Clear-text databases are stored in MySQL version 5.5.

The experiments reported in this section use databases derived from the ClueWeb Collection [18] or synthetically generated by an engine trained on US-census data. The key attributes of these databases and derived encrypted indices are summarized in Table I. The CW-* databases were extracted from the ClueWeb Collection while the LL-* databases emulate the US-census data. Both database families contain atomic type and text columns. The ClueWeb databases were encrypted for a multi-client setting supporting conjunctions (OXT) [3] and the census database where processed for single keyword search (SKS), also in multi-client settings [13] (see Section V-D).

As already mentioned, EDB generation is the dominant factor determining the practical scalability of all our constructions. The two plots called CW (PH) and CW (2L) in Figure 5 show how long it takes to generate the EDBs corresponding to the four CW-* databases when using the Π_{pack} and $\Pi_{2\text{lev}}$ prototypes, respectively.

The results clearly show the $\Pi_{2\text{lev}}$ construction outperforming the Π_{pack} one. The Π_{pack} prototype is consistently slower because its database access patterns are more expensive than those of the $\Pi_{2\text{lev}}$ prototype. For larger datasets, the performance of the Π_{pack} prototype collapses as soon as its RAM requirements, which are proportional with the database size, approach the available RAM. The $\Pi_{2\text{lev}}$ prototype does not exhibit a similar pattern because its RAM requirements are roughly proportional with the size of the database columns currently being processed.

In separate experiments with the $\Pi_{2\text{lev}}$ prototype, preprocessing for the LL-* family of databases also proved to scale linearly, with roughly a rate of $3\mu\text{s}$ per (w, id) pair for the largest database and $8.9\mu\text{s}$ per (w, id) pair for the smallest one. This translates to roughly 92 hours for biggest LL-MC-SKS-4 database as shown in Figure 6. This compares very favorably

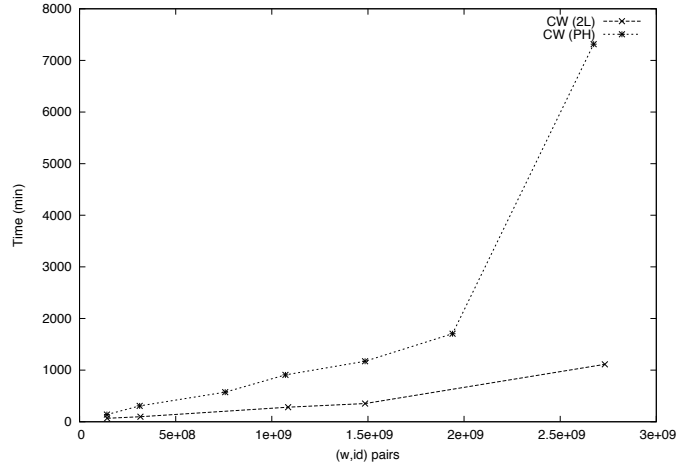


Fig. 5. ClueWeb09 Pre-processing: Scaling Database Size

to the experimental results published by Kamara et al [15], who report a cost of approximately $35\mu\text{s}$ per (w, id) pair on a computing platform with similar capabilities.

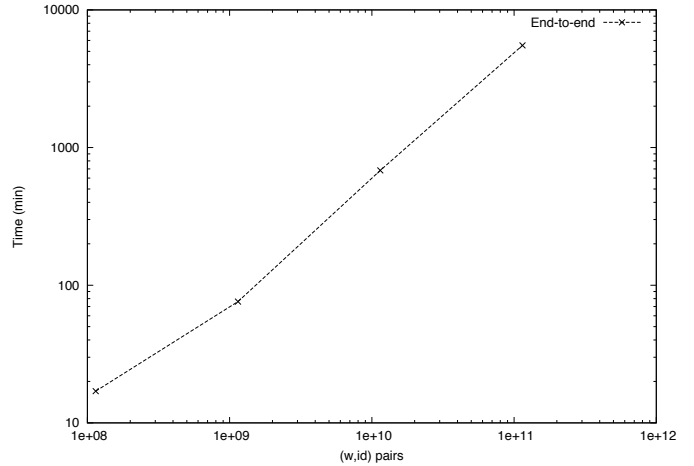


Fig. 6. LL SKS Pre-processing: Scaling Database Size

Measurements on query performance are shown in Figure 7 for queries CW-* databases with varying result set sizes and for both constructions. The graph demonstrates even more dramatic improvements for query processing compared to preprocessing as the $\Pi_{2\text{lev}}$ construction outperforms the Π_{pack} one by almost two orders of magnitude for queries returning 1% of the database. Experiments with the Π_{pack} prototype returning 13% of the database are much slower and they were not included in the figure to improve the visibility of the existing curves. Experiments with OXT demonstrate similar performance gains on conjunction queries. This illustrates that even though OXT performance seemingly is dominated by exponentiation costs (see [3] for the details), optimizing disk-resident data-structures are crucial for good performance due to the high I/O latency costs.⁵

⁵Using highly optimized (common-base) NIST P-224 exponentiation and multi-threading, we can achieve almost 500,000 exponentiation/sec on the mentioned test bed. The storage system provided only, depending on block size, 500-1,500 random I/O requests/sec and single request latencies is around 10ms.

Figure 8 shows the execution times of two queries returning a constant, i.e., independent of the size of the input dataset, result set of 10 and 10,000 record ids, respectively. The gap between the two lines corresponding to the Π_{pack} prototype is much larger than the gap between the lines for corresponding to the $\Pi_{2\text{lev}}$ prototype. The difference between the disk layouts of the two constructs help explain the difference. To retrieve the 10 or the 10,000 ids, the Π_{pack} prototype needs to access one or one thousand hash table buckets, respectively, which means it will issue one thousand times more disk access requests for the second query. In contrast, for the same two queries, the $\Pi_{2\text{lev}}$ prototype needs to access one dictionary entry in both cases and one or eleven *tuple blocks*, which means it will issue only six times more disk access requests for the second query. Π_{pack} hash table buckets and the $\Pi_{2\text{lev}}$ dictionary buckets and *tuple blocks* are all 64KB but tuple groups store only ten tuples in this Π_{pack} prototype while *tuple blocks* store a little under one thousand tuples. Note that since Π_{pack} and $\Pi_{2\text{lev}}$ experiments are using SKS and OXT, respectively, the gap between the 2L and PH plots for experiments returning 10 tuples is explained by the initial computational overhead of OXT.

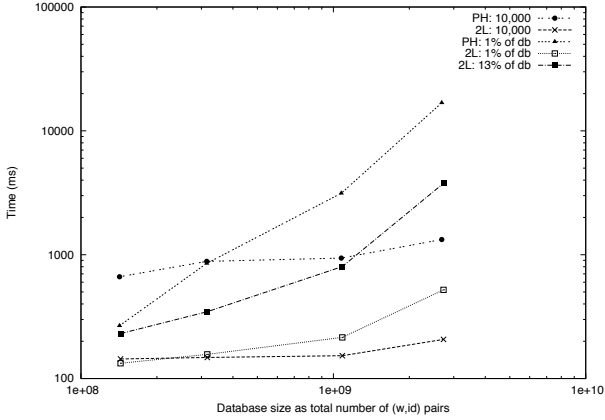


Fig. 7. Clueweb09 SKS Query: Scaling Database Size, comparing Π_{pack} vs $\Pi_{2\text{lev}}$ for various result set sizes.

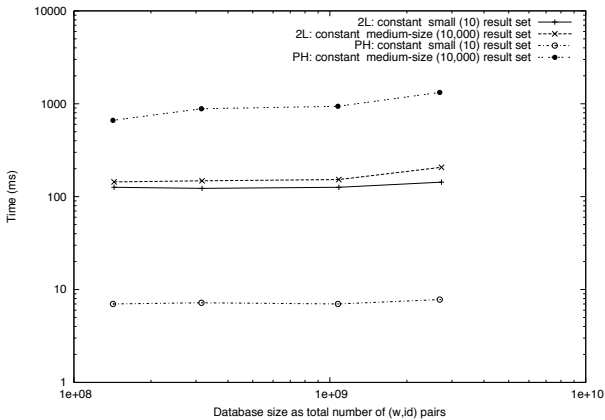


Fig. 8. Clueweb09 SKS Query: Scaling Database Size, comparing Π_{pack} vs $\Pi_{2\text{lev}}$ for constant (10 and 10,000) result set sizes.

Lastly, to illustrate how space efficient the $\Pi_{2\text{lev}}$ construction is, we achieve load-factors of 58.5% for the bucket hash dictionary, 91.9% for the much larger array and 91.6% overall

for our largest LL-MC-SKS-4 database. The load-factor of the array A is less than 100% because although all its entries are used for *tuple blocks*, some of these *tuple blocks* store pointers or are only partially filled with tuples.

F. Comparison with Prior Implementations

The only previous work that stores the encrypted index on external storage is [3], which uses a construction similar to SSE-2 in [7] but adapted to external storage. It corresponds, roughly, to our Π_{pack} prototype discussed in Section V-A. The other previous works assume that the encrypted index is stored in main memory and that access to the index data structure is uniform (i.e., constant cost/latency). None of these constructions admit straightforward extensions to ‘block device’-resident data structures that can be practical. This is particularly the case for constructions using linked lists, such as SSE-1 in [7] or its extension to dynamic SSE in [15].

Recent work by Kamara et. al in [14] discusses an extension of their main memory index to a storage-resident B-tree. This system suffers from using a large index (worst-case quadratic) and their achieved CPU parallelism does not automatically translate to efficient I/O parallelism given the different characteristics of storage sub-systems. The work of [14] does not measure implementation performance and it does not discuss how it would address the I/O challenges faced when building a large scalable system. In contrast, we identify parallel I/O access from the outset as one of the most important performance requirements for scaling to large encrypted indexes. In addition, we also achieve excellent CPU parallelism during search because we parallelize our application-level implementation and because a large fraction of the I/O code path is run in parallel kernel threads. We also validate our approach with experimental results, paramount given the intricacies of storage sub-systems. Finally, our construction does not require a fixed keyword set and is asymptotically faster by $\log n$ than the tree construction in [14], as we use hash instead of tree indexing.

VI. CONCLUSIONS

The tension between security and performance requirements for SSE systems pose non-trivial challenges for both the cryptographic design and the data structures needed to support this design, as well as for their implementation. Leakage minimization calls for randomization of data locations in the encrypted database, EDB, in order to obscure any relations in the original clear-text database. This results in the need to randomize access to EDB elements even when these elements are logically correlated (e.g., the set of documents containing a given keyword). This random-order access is affordable for RAM-resident EDB but becomes prohibitive for disk-resident ones; on the other hand, restricting an implementation to a RAM-resident EDB means limiting the database sizes one can support. Thus, much of the work reported here, both at the abstract data structure level and the specifics of their implementation, are driven by the need to bridge over this security-performance conundrum, and is intended to find a balance between randomized ordering of data, locality of access and parallel processing. In particular, our two-level scheme seems to achieve a desirable trade-off between these competing requirements.

As a result we are able to demonstrate the practicality of search on encrypted data for very large datasets (at terabyte-scale and 10s of billions of record-keyword pairs) and with strong security guarantees. Moreover, our implementation experience shows that even complex queries, as those supported in the work of [3], that go well beyond the single-keyword search capability of traditional SSE schemes, can be supported in practical ways for these very large databases. The same is true for the complex operational settings of [13] that support delegation and authorization of queries to multiple clients as well as providing query privacy from the data owner.

Regarding the security of our schemes, it is important to highlight that while the leakage of our constructions compares well with related work, there is non-trivial information being disclosed to the server about result sets (e.g., the size of these sets and their intersections). When the server has a priori information on the database and queries, various statistical attacks are plausible as illustrated, for instance, in [12]. To mitigate such attacks one can apply various generic masking counter-measures such as padding the sets $DB(w)$ with dummy elements or batching multiple updates to obfuscate update patterns. Hopefully, future work will shed more light on the best ways to design such masking techniques. In particular, one confronts the hard problem of how the syntactically-defined leakage can be captured in a semantic way such that for real world data sets and query distributions one can decide how much and what type of masking approaches are effective.

As mentioned in the introduction, an attractive alternative to achieve more secure solutions to the SSE problem is the use of Oblivious RAM (ORAM) for which we have seen tremendous progress recently in terms of practical performance. However, nobody to our knowledge has yet systematically assessed on how to implement leakage-free search algorithms on top of ORAM servers. Even if we would tolerate the amount of leakage equivalent to our constructions, it is not clear whether one could achieve a similar level of performance for ORAM when considering critical practical aspects such as parallelism and interleaving of I/O and computation as exploited in our approach. Furthermore, the extensibility of ORAM-based solutions to scenarios such as multi-client poses even further challenges.

ACKNOWLEDGMENT

Supported by the Intelligence Advanced Research Projects Activity (IARPA) via Department of Interior National Business Center (DoI / NBC) contract number D11PC20201. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon. Disclaimer: The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of IARPA, DoI/NBC, or the U.S. Government.

REFERENCES

[1] M. Bellare and P. Rogaway. Random oracles are practical: A paradigm for designing efficient protocols. In V. Ashby, editor, *ACM CCS 93*, pages 62–73, Fairfax, Virginia, USA, Nov. 3–5, 1993. ACM Press. 3

[2] D. Boneh, G. Di Crescenzo, R. Ostrovsky, and G. Persiano. Public key encryption with keyword search. In C. Cachin and J. Camenisch, editors, *EUROCRYPT 2004*, volume 3027 of *LNCS*, pages 506–522, Interlaken, Switzerland, May 2–6, 2004. Springer, Berlin, Germany. 2

[3] D. Cash, S. Jarecki, C. S. Jutla, H. Krawczyk, M.-C. Rosu, and M. Steiner. Highly-scalable searchable symmetric encryption with support for boolean queries. In R. Canetti and J. A. Garay, editors, *CRYPTO 2013, Part I*, volume 8042 of *LNCS*, pages 353–373, Santa Barbara, CA, USA, Aug. 18–22, 2013. Springer, Berlin, Germany. 1, 2, 4, 10, 11, 12, 13, 14, 15, 16

[4] Y.-C. Chang and M. Mitzenmacher. Privacy preserving keyword searches on remote encrypted data. In J. Ioannidis, A. Keromytis, and M. Yung, editors, *ACNS 05*, volume 3531 of *LNCS*, pages 442–455, New York, NY, USA, June 7–10, 2005. Springer, Berlin, Germany. 1

[5] M. Chase and S. Kamara. Structured encryption and controlled disclosure. In M. Abe, editor, *ASIACRYPT 2010*, volume 6477 of *LNCS*, pages 577–594, Singapore, Dec. 5–9, 2010. Springer, Berlin, Germany. 1, 2, 3, 4, 5

[6] B. Chor, O. Goldreich, E. Kushilevitz, and M. Sudan. Private information retrieval. *Journal of the ACM*, 45(6):965–981, 1998. 2

[7] R. Curtmola, J. A. Garay, S. Kamara, and R. Ostrovsky. Searchable symmetric encryption: improved definitions and efficient constructions. In A. Juels, R. N. Wright, and S. Vimercati, editors, *ACM CCS 06*, pages 79–88, Alexandria, Virginia, USA, Oct. 30 – Nov. 3, 2006. ACM Press. 1, 2, 3, 4, 15

[8] M. Dietzfelbinger, M. Mitzenmacher, and M. Rink. Cuckoo hashing with pages. Technical Report abs/1104.5111, arXiv, 2011. 11

[9] E.-J. Goh. Secure indexes. Cryptology ePrint Archive, Report 2003/216, 2003. <http://eprint.iacr.org/>. 1

[10] O. Goldreich and R. Ostrovsky. Software protection and simulation on oblivious RAMs. *Journal of the ACM*, 43(3):431–473, 1996. 2

[11] M. T. Goodrich, M. Mitzenmacher, O. Ohrimenko, and R. Tamassia. Oblivious ram simulation with efficient worst-case access overhead. In *CCSW*, pages 95–100, 2011. 3

[12] M. Islam, M. Kuzu, and M. Kantarcioglu. Access pattern disclosure on searchable encryption: Ramification, attack and mitigation. In *Proceedings of the Symposium on Network and Distributed Systems Security (NDSS 2012)*, San Diego, CA, Feb. 2012. Internet Society. 16

[13] S. Jarecki, C. Jutla, H. Krawczyk, M. C. Rosu, and M. Steiner. Outsourced symmetric private information retrieval. In *ACM CCS 13*, Berlin, Germany, Nov. 4–8, 2013. ACM Press. 2, 4, 10, 13, 14, 16

[14] S. Kamara and C. Papamanthou. Parallel and dynamic searchable symmetric encryption. In A.-R. Sadeghi, editor, *FC 2013*, volume 7859 of *LNCS*, pages 258–274, Okinawa, Japan, Apr. 1–5, 2013. Springer, Berlin, Germany. 1, 2, 3, 10, 15

[15] S. Kamara, C. Papamanthou, and T. Roeder. Dynamic searchable symmetric encryption. In T. Yu, G. Danezis, and V. D. Gligor, editors, *ACM CCS 12*, pages 965–976, Raleigh, NC, USA, Oct. 16–18, 2012. ACM Press. 1, 2, 3, 4, 10, 14, 15

[16] J. Katz and Y. Lindell. *Introduction to Modern Cryptography (Chapman & Hall/Crc Cryptography and Network Security Series)*. Chapman & Hall/CRC, 2007. 3

[17] K. Kurosawa and Y. Ohtaki. UC-secure searchable symmetric encryption. In A. D. Keromytis, editor, *FC 2012*, volume 7397 of *LNCS*, pages 285–298, Kralendijk, Bonaire, Feb. 27 – Mar. 2, 2012. Springer, Berlin, Germany. 1, 3

[18] Lemur Project. ClueWeb09 dataset. <http://lemurproject.org/clueweb09.php/>. 14

[19] D. X. Song, D. Wagner, and A. Perrig. Practical techniques for searches on encrypted data. In *2000 IEEE Symposium on Security and Privacy*, pages 44–55, Oakland, California, USA, May 2000. IEEE Computer Society Press. 1, 2

[20] E. Stefanov, E. Shi, and D. X. Song. Towards practical oblivious RAM. In *NDSS 2012*, San Diego, California, USA, Feb. 5–8, 2012. The Internet Society. 3

[21] P. van Liesdonk, S. Sedhi, J. Doumen, P. H. Hartel, and W. Jonker. Computationally efficient searchable symmetric encryption. In *Proc. Workshop on Secure Data Management (SDM)*, pages 87–100, 2010. 1, 3, 10