

Highly-Scalable Searchable Symmetric Encryption with Support for Boolean Queries *

David Cash[†] Stanislaw Jarecki[‡] Charanjit Jutla[§]
Hugo Krawczyk[¶] Marcel-Cătălin Roşu^{||} Michael Steiner^{**}

Abstract

This work presents the design, analysis and implementation of the first searchable symmetric encryption (SSE) protocol that supports conjunctive search and general Boolean queries on outsourced symmetrically-encrypted data and that scales to very large databases and arbitrarily-structured data including free text search. To date, work in this area has focused mainly on single-keyword search. For the case of conjunctive search, prior dedicated SSE constructions (not using generic technique such as FHE or ORAM) required work linear in the total number of documents in the database and provided good privacy only for structured attribute-value data, rendering these solutions too slow and inflexible for large practical databases.

In contrast, our solution provides a realistic and practical trade-off between performance and privacy by efficiently supporting very large databases at the cost of moderate and well-defined leakage to the outsourced server (leakage is in the form of data access patterns, never as direct exposure of plaintext data or searched values). Our design follows a careful process of trading security for efficiency which are both quantified via rigorous analysis. We present a detailed formal cryptographic analysis of the privacy and security of our protocols and establish precise upper bounds on the allowed leakage. To demonstrate the real-world practicality of our approach, we provide performance results of a prototype applied to several large representative data sets, including encrypted search over the whole English Wikipedia (and beyond).

*This is the full version of the paper with same title appearing at Crypto'2013.

[†]Rutgers U. Email: david.cash@cs.rutgers.edu.

[‡]U. California Irvine. Email: stasio@ics.uci.edu.

[§]IBM Research. Email: csjutla@us.ibm.com

[¶]IBM Research. Email: hugo@ee.technion.ac.il

^{||}IBM Research. Email: marcelrosu@gmail.com

^{**}IBM Research. Email: steiner@acm.org

Contents

1	Introduction	3
2	Definitions and Tools	6
2.1	SSE Syntax and Security Model	6
2.2	T-Sets	7
2.3	T-Sets and Single Keyword Search	9
3	SSE Schemes for Conjunctive Keyword Search	9
3.1	Basic Cross-Tags (BXT) Protocol	12
3.1.1	Choosing the s-term	13
3.2	Oblivious Cross-Tags (OXT) Protocol	13
4	Processing Boolean Queries with OXT	16
5	Security Analysis of OXT	17
5.1	Hardness assumptions	17
5.2	Warm up: Analysis for non-adaptive 2-conjunctions	18
5.3	Analysis for boolean queries	21
6	T-Set Instantiation	22
7	OXT Implementation and Experimental Results	25
8	Related Work	28
9	Conclusions and Research Questions	30
A	Proof of Lemma 4	33
B	Proof of Theorem 5	34
C	Oblivious Cross-Tags Protocol PXT Using Bilinear Pairings	42
C.1	Augmented SXDH Assumption	43
C.2	PXT Security Theorem	44
C.3	The Simulator for SSE protocol PXT	44

1 Introduction

Outsourcing data storage to external servers (“the cloud”) is a major industry trend that offers great benefits to database owners. At the same time, data outsourcing raises confidentiality and privacy concerns. Simple encryption of outsourced data is a hindrance to search capabilities such as the data owner wanting to search a backup or email archive, or query a database via attribute-value pairs. This problem has motivated much research on advanced searchable encryption schemes that enable searching on the encrypted data while protecting the confidentiality of data and queries.

Searchable symmetric encryption (SSE) is a cryptographic primitive addressing encrypted search. To securely store and search a database with an SSE scheme, a client first uses a special encryption algorithm which produces an encrypted version of the database, including encrypted metadata, that is then stored on an external server. Later, the client can interact with the server to carry out a search on the database and obtain the results (this is called the symmetric setting as there is only one writer to the database, the owner, who uses symmetric encryption – the public key variant of the problem has also been extensively studied, see further related work).

An important line of research (e.g., [36, 19, 13, 15, 14, 27]) gives practical constructions of SSE that support searching for documents that contain a *single* specified keyword. In these schemes, the server’s work scales with the size of the result set (independently of the database size) and the leakage to the server is limited to the set of returned (encrypted) documents and a few global parameters of the system, such as total data size and number of documents. While efficient and offering good privacy, all of these SSE schemes are severely limited in their expressiveness during search: A client can *only* specify a single keyword to search on, and then it receives *all* of the documents containing that keyword. In practical settings, like remotely-stored email or large databases, a single-keyword search will often return a large number of documents that the user must then download and filter herself to find the relevant results.

Conjunctive and Boolean Search. To provide a truly practical search capability, a system needs to at least support conjunctive search, namely, given a set of keywords find all documents that contain all these keywords. Clearly, this problem can be reduced to the single-keyword case by performing a search for each individual keyword and then letting the server or client do the intersection between the resultant document sets. This often results in inefficient searches (e.g., half the database size if one of the conjunctive terms is “gender=male”) and significant leakage (e.g., it reveals the set of documents matching each keyword). Yet, this naïve solution is the only known sublinear solution to SSE conjunctive search (other than those using generic techniques such as FHE or ORAM). All other dedicated solutions require server work that is *linear in the size of the database*. Of these solutions, the one that provides the best privacy guarantees is due to Golle et al. [21], with variants presented in later work, e.g., [3, 11]. They show how to build for each conjunctive query a set of tokens that can be tested against *each* document in the database (more precisely, against an encoded version of the document’s keywords) to identify matching documents. These solutions only leak the set of matching documents (and possibly the set of attributes being searched for) but are unsuited for large databases due to the $O(d)$ work incurred by the server, where d is the number of documents or records in the database. This cost is paid for every search regardless of the size of the result set or the number of documents matching each individual conjunctive term. Moreover, these solutions require either $O(d)$ communication and exponentiations between server and client or $O(d)$ costly pairing operations (as well as dedicated cryptographic assumptions). Another serious limitation of this approach is that it works only for structured attribute-value type databases and does not support free text search. In addition, none of the above solutions extend to general Boolean queries.

The challenge of large databases and the challenge of being imperfect. In this work we investigate solutions to conjunctive queries and general Boolean queries that can be practical even for very large databases where linear search is prohibitively expensive. Our application settings include databases that require search over tens of millions documents (and billions of document-keyword pairs), with search based on attribute-value pairs (as in relational databases) and free text - see below for specific numbers used in evaluating our prototype. **To support such scale in a truly practical way one needs to relax absolute privacy and allow for some leakage beyond the result set.**

As an example, compare the case of a conjunction of two highly-frequent keywords whose intersection returns a small number of documents but whose individual terms are very frequent (e.g., search for “name=David AND gender=Female”), with the case of a conjunction that returns the same number of documents but all the individual terms in the conjunction are themselves infrequent. Search complexity in these two cases, *even in the case of plaintext data (hence in any encrypted solution)*, is likely to be different and noticeable to the searching server, except if searches are artificially padded to a full database search hence leading to $O(d)$ complexity¹. Note that even powerful tools, such as ORAM, that can be used to search on encrypted data in smaller-scale databases already incur non-trivial leakage if the search performance is to be sublinear. Indeed, the mere computational cost, in number of ORAM operations, of a given search is sufficient to distinguish between the two cases above (of all high-frequency conjunctive terms vs. all small-frequency terms) unless the searches are padded to the maximal search size, resulting in $O(d)$ search cost. Thus, resorting to weaker security guarantees is a necessity for achieving practical conjunctive search. Not only this presents design challenges but also raises non-trivial theoretical challenges for analyzing and characterizing in a precise way the form and amount of leakage incurred by a solution.

Ideally, we would like to run the search with complexity proportional to the number of matches of the *least frequent term in the conjunction*, which is the standard of plaintext information retrieval algorithms. In addition, the computational efficiency of database processing and of search is of paramount importance in practice. Generic tools such as FHE [18] or ORAM [20] are too costly for very large databases, although they may be used as sub-components of a solution if applied to small data subsets.

Our Contributions. We develop the first non-generic sublinear SSE schemes supporting conjunctive keyword search (and more general Boolean queries, see below) with a non-trivial combination of security and efficiency. The schemes performance scales to very large datasets and arbitrarily-structured data, including free-text search. We attain efficiency by allowing some forms of access-pattern leakage, but with a much better leakage profile than the naïve solution implied by single-keyword SSE, as discussed above. Further, we establish the security of our solution via an explicit and precise leakage profile and a *proof that this is all the leakage* incurred by this solution. Our formal setting follows a simulation-based abstraction that adapts the SSE models of Curtmola et al. [15] and Chase and Kamara [14], and assumes an adaptive adversarial model. The essence of the security notion is that the view of the server (the attacker in this setting) can be efficiently simulated given a precisely-defined *leakage profile* but without access to the actual plaintext data. Such a profile may include leakage on the total size of the database, on access patterns (e.g., the intersection between two sets of results) and on queries (e.g., repetition of queries), but never the direct exposure of plaintext data or searched values. Thus, a protocol proven secure ensures that the server holding the encrypted data and serving the queries does not learn anything about the data and queries other than what can be deduced from the specified leakage². The characterization

¹A costly alternative is to pre-compute all n -term conjunctions in time $O(|W|^n)$.

² See the discussion in Section 9 on “semantic leakage”.

of leakage and the involved proof of security that we present are central technical contributions that complement our protocol design work.

The centerpiece of the protocol design is a “virtual” secure two-party protocol in which the server holds encrypted pointers to documents, the client holds a list of keywords, and the output of the protocol is the set of encrypted pointers that point to documents containing all the client’s keywords. The client is then able to decrypt these pointers and obtain the matching (encrypted) documents but the server cannot carry this decryption nor can it learn the keywords in the client’s query. While this underlying protocol is interactive, the level of performance targeted by our solutions requires avoiding multiple rounds of interaction. We achieve this by a novel approach that pre-computes parts of the protocol messages and stores them in encrypted form at the server. Then, during search, the client sends information to the server that allows to unlock these pre-computed messages without further interaction. Our implementation of this protocol, which we name OXT, uses only DH-type operations over any Diffie-Hellman group which enables the use of the secure and most efficient DH elliptic curves (with additional common-base optimizations).³ The complexity of our search protocols is *independent* of the number of documents in the database. To search for documents containing w_1, \dots, w_n , the search complexity of our scheme scales with the number of documents matching the estimated *least frequent keyword in the conjunction*. We note that while building a search based on term frequency is standard in information retrieval, our solution seems to be the first to exploit this approach in the encrypted setting. This leads not only to good performance but also improves privacy substantially. All our solutions support search on structured data (e.g., attribute-value databases) as well as on free text, and combinations of both.

Boolean queries. Our solution to conjunctive queries extends to answer *any Boolean query*. This includes negations, disjunctions, threshold queries, and more. The subset of such queries that we can answer efficiently includes any expression of the form “ $w_1 \wedge \phi(w_2, \dots, w_m)$ ” (intended to return any document that matches keyword w_1 and in addition satisfies the (unconstrained) formula ϕ on the remaining keywords)⁴. The search complexity is proportional to the number of documents that contain w_1 . Surprisingly, the leakage profile for such complex expressions can be reduced to the leakage incurred by a conjunction with the same terms w_1, w_2, \dots, w_n , hence allowing us to re-use the analysis of the conjunctive case to the much more general boolean setting. Finally, any disjunction of the above forms can also be answered with an additive cost over the disjunction expressions.

Further extensions. In [12] we report on further practical enhancements to our protocols, including support for dynamic databases (i.e., allowing additions, deletions and modification of documents in the database) and increased scalability. Our protocols can also be applied to the *multi-client setting* [14, 25, 26] where a data owner outsources its encrypted data to an external server and enables *other parties* to perform queries on the encrypted data by providing them with search tokens for specific queries. In this case, one considers not only leakage to the server but also leakage to clients beyond the information that their tokens are authorized to disclose. In subsequent work [24] we address issues of authorization in this setting as well as the challenging problem of hiding the queries not only from the server but also from the token provider - see for example IARPA’s SPAR program [22] and its requirement for supporting private queries on very large databases. See also [31] for an independent, concurrent work in the latter setting from which a solution to the SSE problem can also be extracted. Finally, in ongoing work, we are extending the set of supported

³ We also present a scheme (BXT in Section 3.1) that only uses symmetric-key operations but provides less privacy, and a pairings-based scheme (PXT in Section C) that optimizes communication at the expense of more computation.

⁴ An example of such query on an email repository is: Search for messages with Alice as Recipient, not sent by Bob, and containing at least two of the words {searchable, symmetric, encryption}.

queries with range queries, substring matching, and more.

Implementation. To show the practical viability of our solution we prototyped OXT and ran experiments with three data sets: a 100,000 record relational database synthesized from census data; the Enron email data set [16] with more than 1.5 million documents (email messages and attachments) where all words, including attachments and envelope information, have been indexed; and the ClueWeb09 [29] collection of crawled web-pages from which we extracted several databases of increasing size with the largest one consisting of 13 million documents (0.4TB of HTML files). Approximately one third of the latter database is a full snapshot of the English Wikipedia. The results of these tests show not only the suitability of our conjunction protocols for data sets of medium size (such as the Enron one) but demonstrate the scalability of these solutions to much larger databases (we target databases of one or two orders of magnitude larger). Existing solutions that are linear in the number of documents would be mostly impractical even for the Enron dataset. Refer to Section 7 for more information on implementation and performance. More advanced results are reported in [12], including ways to handle data updates in dynamic databases as well as supporting DBs of up to two orders of magnitude above the experiments reported here while preserving the full Boolean query functionality.

Other related work and research questions. See Section 8 for more discussion on related work and Section 9 for several interesting research questions arising from our work.

2 Definitions and Tools

Notation. We write $[n]$ for the set $\{1, \dots, n\}$. For a vector \mathbf{v} we write $|\mathbf{v}|$ for the dimension (length) of \mathbf{v} and for $i \in [|\mathbf{v}|]$ we write $\mathbf{v}[i]$ for the i -th component of \mathbf{v} . All algorithms (including adversaries) are assumed to be randomized polynomial-time unless otherwise specified. If A is an algorithm, then $y \leftarrow A(x)$ means that the y is the output of A when run on input x . If A is randomized then y is a random variable. For sets X, Y we write $\text{Fun}(X, Y)$ for the set of all functions from X to Y , and $\text{Perm}(X)$ for the set of all permutations on X .

2.1 SSE Syntax and Security Model

Searchable symmetric encryption. A database is composed of a collection of d documents, each comprised of a set of keywords W_i (we use “documents” generically; they can represent text documents, records in a relational database - in which case keyword are represented as attribute-value pairs, a combination of both, etc.). The output from the SSE protocol for a given search query are indices (or identifiers) ind corresponding to the documents that satisfy the query. A client program can then use these indices to retrieve the encrypted documents and decrypt them. This definition allows to **decouple** the storage of payloads (which can be done in a variety of ways, with varying types of leakage) **from** the storage of metadata that is the focus of our protocols.

SSE scheme syntax and correctness. Let λ be the security parameter. We will take identifiers and keywords to be bit strings. A database $\text{DB} = (\text{ind}_i, \mathbf{W}_i)_{i=1}^d$ is represented as a list of identifier/keyword-set pairs, where $\text{ind}_i \in \{0, 1\}^\lambda$ and $W_i \subseteq \{0, 1\}^*$. We will always write $W = \bigcup_{i=1}^d W_i$ (we think of the ind values as identifiers that can be revealed to the outsourced server, e.g., a randomization of the original document identifiers; in particular these are the identifiers that will be used to retrieve query-matching documents). A *query* $\psi(\bar{w})$ is specified by a tuple of keywords $\bar{w} \in W^*$ and a boolean formula ψ on \bar{w} . We write $\text{DB}(\psi(\bar{w}))$ for the set of identifiers of documents that “satisfy” $\psi(\bar{w})$. Formally, this means that $\text{ind}_i \in \text{DB}(\psi(\bar{w}))$ iff the formula $\psi(\bar{w})$



evaluates to true when we replace each keyword w_i with true or false depending on if $w_i \in W_i$ or not. Below we let d denote the number of documents in DB, $m = |W|$ and $N = \sum_{w \in W} |\text{DB}(w)|$.

A *searchable symmetric encryption (SSE) scheme* Π consists of an algorithm **EDBSetup** and a protocol **Search** between the client and server, all fitting the following syntax. **EDBSetup** 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 $\psi(\bar{w})$ and the server takes as input EDB. At the end of the protocol the client outputs a set of identifiers and the server has no output. We say that an SSE scheme is *correct* if for all inputs DB and queries $\psi(\bar{w})$ for $\bar{w} \in W^*$, if $(K, \text{EDB}) \xleftarrow{\$} \text{EDBSetup}(\text{DB})$, after running **Search** with client input $(K, \psi(\bar{w}))$ and server input EDB, the client outputs the set of indices $\text{DB}(\psi(\bar{w}))$. We consider a computational relaxation of this notion, expressed via the following game. For an adversary A and an SSE scheme Σ , we define the game $\text{Cor}_A^\Pi(\lambda)$, which lets A choose DB, generates $(K, \text{EDB}) \leftarrow \text{EDBSetup}(\text{DB})$, gives EDB to A , which adaptively chooses queries $\psi(\bar{w})$, for each of which the game runs the **Search** protocol with client input $(K, \psi(\bar{w}))$ and server input EDB. If in any execution the client outputs something other than $\text{DB}(\psi(\bar{w}))$, the game outputs 1, and otherwise it outputs 0. We say that a scheme Π is *computationally correct* if for all efficient adversaries A , $\Pr[\text{Cor}_A^\Pi(\lambda) = 1] \leq \text{neg}(\lambda)$.

Security of SSE. We recall the semantic security definitions from [15, 14]. The definition is parametrized by a *leakage function* \mathcal{L} , which describes what an adversary (the server) is allowed to learn about the database and queries when interacting with a secure scheme. Formally, security says that the server's view during an adaptive attack (where the server selects the database and queries) can be simulated given only the output of \mathcal{L} .

Definition 1 Let $\Pi = (\text{EDBSetup}, \text{Search})$ be an SSE scheme and let \mathcal{L} be a stateful algorithm. For algorithms A and S , we define experiments (algorithms) $\text{Real}_A^\Pi(\lambda)$ and $\text{Ideal}_{A,S}^\Pi(\lambda)$ as follows:

Real_A^Π(λ) : $A(1^\lambda)$ chooses DB. The experiment then runs $(K, \text{EDB}) \leftarrow \text{EDBSetup}(\text{DB})$ and gives EDB to A . Then A repeatedly chooses a query q . To respond, the game runs the **Search** protocol with client input (K, q) and server input EDB and gives the transcript and client output to A . Eventually A returns a bit that the game uses as its own output.

Ideal_{A,S}^Π(λ) : The game initializes a counter $i = 0$ and an empty list \mathbf{q} . $A(1^\lambda)$ chooses DB. The experiment runs $\text{EDB} \leftarrow S(\mathcal{L}(\text{DB}))$ and gives EDB to A . Then A repeatedly chooses a query q . To respond, the game records this as $\mathbf{q}[i]$, increments i , and gives to A the output of $S(\mathcal{L}(\text{DB}, \mathbf{q}))$. (Note that here, \mathbf{q} consists of all previous queries in addition to the latest query issued by A .) Eventually A returns a bit that the game uses as its own output.

We say that Π is \mathcal{L} -semantically-secure against adaptive attacks if for all adversaries A there exists an algorithm S such that $\Pr[\text{Real}_A^\Pi(\lambda) = 1] - \Pr[\text{Ideal}_{A,S}^\Pi(\lambda) = 1] \leq \text{neg}(\lambda)$.

We note that in the security analysis of our SSE schemes we include the client's output, the set of indices $\text{DB}(\psi(\bar{w}))$, in the adversary's view in the real game, to model the fact that these ind's will be used for retrieval of encrypted document payloads. In Appendix B we include a version of this security notion to *non-adaptive* adversaries.

2.2 T-Sets

We present a definition of syntax and security for a new primitive that we call a *tuple set*, or *T-set*. Intuitively, a T-set allows one to associate a list of fixed-sized data tuples with each keyword in

the database, and later issue keyword-related tokens to retrieve these lists. We will use it in our protocols as an “expanded inverted index”. Indeed, prior single-keyword SSE schemes, e.g. [15, 14], can be seen as giving a specific T-set instantiation and using it as an inverted index to enable search – see Section 2.3. In our SSE schemes for conjunctive keyword search, we will use a T-set to store more data than a simple inverted index, and we will also compose it with other data structures. The abstract definition of a T-set will allow us to select an instantiation that provides the best performance for the size of the data being stored. One specific T-set instantiation is shown in Section 6.

T-Set syntax and correctness. Formally, a T-set instantiation $\Sigma = (\text{TSetSetup}, \text{TSetGetTag}, \text{TSetRetrieve})$ will consist of three algorithms with the following syntax: TSetSetup will take as input an array \mathbf{T} of lists of equal-length bit strings indexed by the elements of \mathbf{W} . In other words, for any function $n(\lambda)$ of the security parameter λ , for each $w \in \mathbf{W}$, $\mathbf{T}[w]$ is a list $\mathbf{t} = (s_1, \dots, s_{T_w})$ of strings s.t. each s_i is of length $n(\lambda)$, but $T_w = |\mathbf{T}[w]|$ may vary with w . (In our applications of T-set, $\mathbf{T}[w]$ will contain one tuple per each DB document which matches w , i.e. $T_w = |\text{DB}(w)|$.) The TSetSetup procedure outputs a pair (TSet, K_T) . TSetGetTag takes as input the key K_T and a keyword w and outputs stag . TSetRetrieve takes the TSet and an stag as input, and returns a list of strings. We say that Σ is *correct* if for all \mathbf{W} , \mathbf{T} , and any $w \in \mathbf{W}$, $\text{TSetRetrieve}(\text{TSet}, \text{stag}) = \mathbf{T}[w]$ when $(\text{TSet}, K_T) \leftarrow \text{TSetSetup}(\mathbf{T})$ and $\text{stag} \leftarrow \text{TSetGetTag}(K_T, w)$. Intuitively, \mathbf{T} holds lists of tuples associated with keywords and correctness guarantees that the TSetRetrieve algorithm returns the data associated with the given keyword. As in the case of SSE, we consider a computational relaxation of this perfect correctness notion: For adversary A and a T-set instantiation Σ , we define a game $\text{AdvCor}_A^\Pi(\lambda)$, where A chooses \mathbf{T} , the game generates $(\text{TSet}, K_T) \leftarrow \text{TSetSetup}(\mathbf{T})$, gives TSet to A , which adaptively chooses keywords w , for each of which the game generates $\text{stag} \leftarrow \text{TSetGetTag}(K_T, w)$ and $\mathbf{t}_w \leftarrow \text{TSetRetrieve}(\text{TSet}, \text{stag})$. The game outputs 1 if for any w we have $\mathbf{t}_w \neq \mathbf{T}[w]$, and 0 otherwise. We say that a T-set instantiation Π is *computationally correct* if for all efficient adversaries A , $\text{AdvCor}_A^\Pi(\lambda)$, defined as the probability that the above game outputs 1, is a negligible function of λ .

T-Set security and instantiation. The security goal of a T-set instantiation is to hide as much as possible about the tuples in \mathbf{T} and the keywords these tuples are associated to, except for vectors $\mathbf{T}[w_1], \mathbf{T}[w_2], \dots$ of tuples revealed by the client’s queried keywords w_1, w_2, \dots . (For the purpose of T-set instantiation we equate client’s query with a single keyword.) Since the list of tuples associated to searched keywords can be seen as information provided to the server, this information is provided to the simulator in the security definition below.

We parametrize the T-set security definition with a leakage function \mathcal{L}_T that describes what else the adversary is allowed to learn by looking at the TSet and stag values. For most instantiations this leakage will reveal something about the structure of \mathbf{T} , and consequently also the structure of DB . For example, an instantiation could reveal the size of \mathbf{T} , which is the number of keywords in DB , or the length of each list $\mathbf{T}[w]$ in \mathbf{T} , which reveals the number of occurrences of each keyword w in DB . Our more careful instantiation given in Section 6 can be shown to leak significantly less, namely only $N = \sum_{w \in \mathbf{W}} |\mathbf{T}[w]|$, the total number of keyword occurrences in DB .

Definition 2 Let $\Sigma = (\text{TSetSetup}, \text{TSetGetTag}, \text{TSetRetrieve})$ be a T-set instantiation, and let A, S be an adversary and a simulator, and let \mathcal{L}_T be a stateful algorithm. We define two games, Real_A^Σ and Ideal_A^Σ as follows.

$\text{Real}_A^\Sigma(\lambda) : A(1^\lambda)$ outputs \mathbf{W}, \mathbf{T} with the above syntax. The game computes $(\text{TSet}, K_T) \leftarrow \text{TSetSetup}(\mathbf{T})$ and gives TSet to A . Then A repeatedly issues queries $q \in \mathbf{W}$, and for each q the game gives $\text{stag} \leftarrow \text{TSetGetTag}(K_T, q)$ to A . Eventually A outputs a bit which the game uses as its output.

Ideal $_{A,S}^{\Sigma}(\lambda)$: The game initializes a counter $i = 0$ and an empty list \mathbf{q} . $A(1^{\lambda})$ outputs \mathbf{W}, \mathbf{T} as above. The game runs $\mathbf{TSet} \leftarrow S(\mathcal{L}_T(\mathbf{T}))$ and gives \mathbf{TSet} to A . Then A repeatedly issues queries $q \in \mathbf{W}$, and for each q the game stores q in $\mathbf{q}[i]$, increments i , and gives to A the output of $S(\mathcal{L}_T(\mathbf{T}, \mathbf{q}), \mathbf{T}[q])$. Eventually A outputs a bit which the game uses as its output.

We say that Σ is a \mathcal{L}_T -adaptively-secure T-set instantiation if for all adversaries A there exists an algorithm S such that $\Pr[\mathbf{Real}_A^{\Sigma}(\lambda) = 1] - \Pr[\mathbf{Ideal}_{A,S}^{\Sigma}(\lambda) = 1] \leq \text{neg}(\lambda)$.

A non-adaptive version of this definition is a straightforward modification of the above game where the adversary provides all of the queries at the start of the game.

2.3 T-Sets and Single Keyword Search

Here we show how a T-set can be used as an “secure inverted index” to build an SSE scheme for single-keyword search. The ideas in this construction will be the basis for our conjunctive search SSE schemes later, and it essentially abstracts prior constructions [15, 14]. The details of the scheme, called **SKS**, are given in Figure 1. It uses as subroutines a PRF $F : \{0, 1\}^{\lambda} \times \{0, 1\}^{\lambda} \rightarrow \{0, 1\}^{\lambda}$, and a CPA secure symmetric encryption scheme (Enc, Dec) that has λ -bit keys. When combined with the specific adaptively-secure T-set construction shown in Section 6, the **SKS** protocol results in a single-keyword SSE scheme that achieves the best storage and search performance among existing adaptive SSE protocols. Whereas our **EDB** stores $c * N$ short ciphertexts for a small constant $c \leq 3$, where $N = \sum_{w \in \mathbf{W}} |\text{DB}(w)|$, i.e. the total number of keyword occurrences in DB , the adaptive SSE scheme of [15] needs $d * \max$ storage where d is the number of documents and \max is the maximum number of keywords in any document, while the adaptive SSE scheme for keyword search of [14] needs $m * \max'$ storage where m is the number of keywords and \max' is the maximum number of documents matching any keyword. Consequently our T-set implementation improves over adaptive single-keyword SSE of [15] by factor equal to the proportion between maximum number of keywords in any document to the average number of keywords per document, and it improves over adaptive single-keyword SSE of [14] by factor equal to the proportion between maximum number of documents matching any keyword to the average number of documents per keyword.⁵ Finally, the storage complexity of [27] is similar to ours, but their search procedure cannot be parallelized on the server because they represent a T-set as a linked list, and their reliance on the ROM model for security seems harder to do away with. (On the other hand, the scheme of [27] enables updates without growth in the data structure size and search time.)

Recent work [12] improves on our Tset construction offering a two-level scheme that optimizes both disk access and memory load factors scaling to almost two orders of magnitude larger DBs than in the present work (while also providing support for dynamic databases).

3 SSE Schemes for Conjunctive Keyword Search

Existing SSE schemes for conjunctive queries ([21] and subsequent work) work by encoding each document individually and then processing a search by testing *each* encoded document against a set of tokens. Thus the server’s work grows linearly with the number of documents, which is infeasible

⁵As for communication costs, the adaptive scheme of [15] has $O(d)$ communication as stated, but it can be easily reduced to $O(|\text{DB}(w)|)$, the size of the result set. The scheme of [14] has $O(\max')$ communication, i.e. $O(|\text{DB}(w)|)$ for the worst case w . The schemes of [15, 14] do not rely on ROM, but our T-set implementation can also avoid ROM if instead of a single value **strap** the client streams to the server consecutive values $F(\mathbf{strap}, 1), F(\mathbf{strap}, 2), \dots$ until receiving server’s signal to stop.

EDBSetup(DB)

- Select key K_S for PRF F , and parse DB as $(\text{ind}_i, W_i)_{i=1}^d$.
- Initialize \mathbf{T} to an empty array indexed by keywords from $W = \cup_{i=1}^d W_i$.
- For each $w \in W$, build the tuple list $\mathbf{T}[w]$ as follows:
 - Initialize \mathbf{t} to be an empty list, and set $K_e \leftarrow F(K_S, w)$.
 - For all $\text{ind} \in \text{DB}(w)$ in random order: $\mathbf{e} \xleftarrow{\$} \text{Enc}(K_e, \text{ind})$; append \mathbf{e} to \mathbf{t} .
 - Set $\mathbf{T}[w] \leftarrow \mathbf{t}$.
- $(\text{TSet}, K_T) \leftarrow \text{TSetSetup}(\mathbf{T})$.
- Output the key (K_S, K_T) and $\text{EDB} = \text{TSet}$.

Search protocol

- The client takes as input the key (K_S, K_T) and a keyword w to query. It computes $\text{stag} \leftarrow \text{TSetGetTag}(K_T, w)$ and sends stag to the server.
- The server computes $\mathbf{t} \leftarrow \text{TSetRetrieve}(\text{TSet}, \text{stag})$, and sends \mathbf{t} to the client.
- Client sets $K_e \leftarrow F(K_S, w)$; for each \mathbf{e} in \mathbf{t} , it computes $\text{ind} \leftarrow \text{Dec}(K_e, \mathbf{e})$ and outputs ind .

Figure 1: SKS: SINGLE-KEYWORD SSE SCHEME

for large databases. In addition, these schemes only work for attribute-value type databases (where documents contain a single value per attribute) but not for unstructured data, e.g., they cannot search text documents.

Here we develop the first sub-linear conjunctive-search solutions for arbitrarily-structured data, including free text. In particular, when querying for the documents that match all keywords w_1, \dots, w_n , our search protocol scales with the size of the (estimated) *smallest* $\text{DB}(w_i)$ set among all the conjunctive terms w_i .

The naïve solution. To motivate our solutions we start by describing a straightforward extension of the single-keyword case (protocol SKS from Figure 1) to support conjunctive keyword searching. On input a conjunctive query $\bar{w} = (w_1, \dots, w_n)$, the client and server run the search protocol from SKS independently for each term w_i in \bar{w} with the following modifications. Instead of returning the lists \mathbf{t} to the client, the server receives K_{e_i} , $i = 1, \dots, n$, from the client and decrypts the \mathbf{e} values to obtain a set of ind 's for each w_i . Then, the server returns to client the ind values in the intersection of all these sets. The search complexity of this solution is proportional to $\sum_{i=1}^n |\text{DB}(w_i)|$ which improves, in general, on solutions whose complexity is linear in the number of documents in the whole database. However, this advantage is reduced for queries where one of the terms is a very high-frequency word (e.g., in a relational database of personal records, one may have a keyword $w = (\text{gender}, \text{male})$ as a conjunctive term, thus resulting in a search of, say, half the documents in the database). In addition, this solution incurs excessive leakage to the server who learns the complete sets of indices ind for each term in a conjunction.

Our goal is to reduce both computation and leakage in the protocol by tying those to the less frequent

terms in the conjunctions (i.e., terms w with small sets $\text{DB}(w)$).

Text in red indicates changes from SKS, Figure 1

EDBSetup(DB)

- Select keys K_S and K_X for PRF F ; parse DB as $(\text{ind}_i, W_i)_{i=1}^d$.
- Initialize \mathbf{T} to an empty array indexed by keywords from W .
- Initialize XSet to an empty set.
- For each $w \in W$, build the tuple list $\mathbf{T}[w]$ and XSet elements as follows:
 - Initialize \mathbf{t} to be an empty list, and set $K_e \leftarrow F(K_S, w)$.
 - Compute $\mathbf{xtrap} \leftarrow F(K_X, w)$
 - For all ind in $\text{DB}(w)$ in random order:
 - * Compute $\mathbf{e} \leftarrow \text{Enc}(K_e, \text{ind})$ and append \mathbf{e} to \mathbf{t} .
 - * $\mathbf{xtag} \leftarrow f(\mathbf{xtrap}, \text{ind})$ and add \mathbf{xtag} to XSet.
 - $\mathbf{T}[w] \leftarrow \mathbf{t}$.
- $(\text{TSet}, K_T) \leftarrow \text{TSetSetup}(\mathbf{T})$.
- Output the key (K_S, K_X, K_T) and $\text{EDB} = (\text{TSet}, \text{XSet})$.

Search protocol

- The client takes as input the key (K_S, K_X, K_T) and keywords w_1, \dots, w_n to query. It computes its messages as
 - $K_e \leftarrow F(K_S, w_1)$, $\text{stag} \leftarrow \text{TSetGetTag}(K_T, w_1)$, then
 - For each $i = 2, \dots, n$, it sets $\mathbf{xtrap}_i \leftarrow F(K_X, w_i)$.
 - It sends $(\text{stag}, K_e, \mathbf{xtrap}_2, \dots, \mathbf{xtrap}_n)$ to the server.
- The server has input $(\text{TSet}, \text{XSet})$. It responds as follows.
 - It sets $\mathbf{t} \leftarrow \text{TSetRetrieve}(\text{TSet}, \text{stag})$.
 - For each ciphertext \mathbf{e} in \mathbf{t} , it computes
 - * $\text{ind} \leftarrow \text{Dec}(K_e, \mathbf{e})$
 - * If $f(\mathbf{xtrap}_i, \text{ind}) \in \text{XSet}$ for all $i = 2, \dots, n$, it sends ind to the client.
- The client outputs all of the received inds.

Figure 2: BXT: BASIC CROSS-TAGS PROTOCOL

3.1 Basic Cross-Tags (BXT) Protocol

To achieve the above goal we take the following approach that serves as the basis for our main SSE-conjunctions scheme OXT presented in the next subsection. Here we exemplify the approach via a simplified protocol, BXT. Assume (see Section 3.1.1) that the client, given $\bar{w} = (w_1, \dots, w_n)$, can choose a term w_i with a relatively small $\text{DB}(w_i)$ set among w_1, \dots, w_n ; for simplicity assume this is w_1 . The parties could run an instance of the SKS search protocol for the keyword w_1 after which the client gets all documents matching w_1 and locally searches for the remaining conjunctive terms. This is obviously inefficient as it may require retrieving many more documents than actually needed. The idea of BXT is indeed to use SKS for the server to retrieve $\text{TSet}(w_1)$ but then perform the intersection with the terms w_2, \dots, w_n at the server who will only return the documents matching the full conjunction. We achieve this by augmenting SKS as follows.

During $\text{EDBSetup}(\text{DB})$, in addition to TSet , a set data structure XSet is built by adding to it elements xtag computed as follows. For each $w \in W$, a value $\text{xtrap} = F(K_X, w)$ is computed where K_X is a PRF key chosen for this purpose; then for each $\text{ind} \in \text{DB}(w)$ a value $\text{xtag} = f(\text{xtrap}, \text{ind})$ is computed and added to XSet where f is an unpredictable function of its inputs (e.g., f can be a PRF used with xtrap as the key and ind as input). The Search protocol for a conjunction (w_1, \dots, w_n) , chooses the estimated least frequent keyword, say w_1 , and sets, as in SKS, $K_e \leftarrow F(K_S, w_1)$, $\text{stag} \leftarrow \text{TSetGetTag}(K_T, w_1)$. Then, for each $i = 2, \dots, n$, it sets $\text{xtrap}_i \leftarrow F(K_X, w_i)$ and sends $(K_e, \text{stag}, \text{xtrap}_2, \dots, \text{xtrap}_n)$ to the server. The server uses stag to retrieve $\mathbf{t} = \text{TSetRetrieve}(\text{TSet}, \text{stag})$. Then, for each ciphertext \mathbf{e} in \mathbf{t} , it decrypts $\text{ind} = \text{Dec}(K_e, \mathbf{e})$ and if $f(\text{xtrap}_i, \text{ind}) \in \text{XSet}$ for all $i = 2, \dots, n$, it sends ind to the client.⁶

Correctness of the BXT protocol is easy to verify. Just note that a document indexed by ind includes a word w represented by stag if and only if $\text{xtag} = f(\text{xtrap}, \text{ind}) \in \text{XSet}$. Regarding implementation of XSet , it can use any set representation that is history-independent, namely, it is independent of the order in which the elements of the set were inserted. For TSet security and implementation see Section 2.

Terminology (s-terms and x-terms): We will refer to the conjunctive term chosen as the estimated least frequent term among the query terms as the *s-term* (‘s’ for SKS or “small”) and refer to other terms in the conjunction as *x-terms* (‘x’ for “cross”); this is the reason for the ‘s’ and ‘x’ in names such as stag , xtag , stag , xtrap , etc.

The server’s work in BXT scales with $n \cdot |\text{DB}(w_1)|$, where w_1 is the conjunction’s s-term. This represents a major improvement over existing solutions which are linear in $|\text{DB}|$ and also a significant improvement over the naïve solution whenever there is a term with relatively small set $\text{DB}(w_1)$ that can be identified by the client, which is usually the case as discussed in Section 3.1.1. Communication is optimal ($O(n)$ -size token plus the final results set) and computation involves only PRF operations.

Security-wise this protocol *improves substantially* on the above-described naïve solution by leaking only the (small) set of ind ’s for the s-term and not for x-terms. Yet, this solution lets the server learn statistics about x-terms by correlating information from different queries. Specifically, the server can use the value xtrap_i received in one query and check it against any ind found through an s-term of another query. But note that direct intersections between x-terms of different queries are not possible other than via the s-terms (e.g., if two queries (w_1, w_2) and (w'_1, w'_2) are issued, the server can learn the (randomly permuted) results of (w_1, w'_2) and (w'_1, w_2) but not (w_2, w'_2)).

In settings where computation and communications are very constrained BXT may provide for

⁶ While in SKS one can choose to let the server decrypt the ind ’s directly instead of the client, in BXT this is necessary for computing the xtag ’s.

an acceptable privacy-performance balance. In general, however, we would like to improve on the privacy of this solution even if at some performance cost. We do so in the next section with the OXT protocol, so we omit a formal analysis BXT – we note that the security of BXT needs the set of *ind*’s to be unpredictable, a condition not needed in the other protocols.

3.1.1 Choosing the s-term

The performance and privacy of our conjunction protocols improves with “lighter” s-terms, namely, keywords w whose $\text{DB}(w)$ is of small or moderate size. While it is common to have such terms in typical conjunctive queries, our setting raises the question of how can the client, who has limited storage, choose adequate s-terms. In the case of relational databases one can use general statistics about attributes to guide the choice of the s-term (e.g., prefer a last-name term to a first-name term, always avoid gender as the s-term, etc.). In the case of free text the choice of s-term can be guided by term frequency which can be made available, requiring a small state stored at the client or retrieved from the server. Indeed, it is the common case that the number of very frequent words in a document collection is relatively small. This means that the client can keep a small state with sufficient information to choose light s-terms. For example, in the Enron data set [16] used in our testing, only 4% of 1,176,222 distinct words from 1,551,675 documents appear in more than 100 documents. With a state of less than 100 Kbyte a client can keep information, e.g., via a Bloom filter, to differentiate between keywords with less than 100, 1000, 10,000, 100,000, matching documents, respectively (and if one omits the ‘> 100’ category then the storage requirement reduces to less than 25KB). We note that on our experimental platform, any $\text{TSet}(w)$ with less than 1000 tuples is retrieved in less than a tenth of a second, so the level of granularity for the above tables can be tuned depending on the data and setting. A client that cannot afford storing a state of the above size, can retrieve it from the server (where this state is stored encrypted) at the beginning of a search session. And there is always the option to cap the number of searched and/or retrieved documents (after which the client may choose to refine its search). Finally, we observe that in the case of multi-client SSE setting [14] discussed in the introduction, the data owner, who provides tokens to clients, will usually have enough storage to make optimal (or close-to-optimal) choices of s-terms.

3.2 Oblivious Cross-Tags (OXT) Protocol

The BXT scheme is vulnerable to the following simple attack: When the server gets xtrap_i for a query (w_1, \dots, w_n) , it can save it and later use it to learn if any revealed *ind* value is a document with keyword w_i by testing if $f(\text{xtrap}_i, \text{ind}) \in \text{XSet}$. This allows an honest-but-curious server to learn, for example, the number of documents matching each queried s-term with each queried x-term (even for terms in different queries). This attack is possible because BXT reveals the keys that enable the server to compute $f(\text{xtrap}_i, \cdot)$ itself.

One way to mitigate the attack is to have the client evaluate the function for the server instead of sending the key. Namely, the server would send all the encrypted *ind* values that it gets in \mathbf{t} (from the TSet) to the client who will compute the function $f(\text{xtrap}_i, \text{ind})$ and send back the results. However, this fix adds a round of communication with consequent latency, it allows the server to cheat by sending *ind* values from another query’s s-term (from which the server can compute intersections not requested by the client), and is unsuited to the multi-client SSE setting [14] discussed in the introduction (since the client would learn from the *inds* it receives the results of conjunctions it was not authorized for). Note that while the latter two issues are not reflected in our current formal model, avoiding them expands significantly the applicability of OXT.

EDBSetup(DB)

- Select key K_S for PRF F , keys K_X, K_I, K_Z for PRF F_p (with range in Z_p^*), and parse DB as $(\text{ind}_i, W_i)_{i=1}^d$.
- Initialize \mathbf{T} to an empty array indexed by keywords from W .
- Initialize XSet to an empty set.
- For each $w \in W$, build the tuple list $\mathbf{T}[w]$ and XSet elements as follows:
 - Initialize \mathbf{t} to be an empty list, and set $K_e \leftarrow F(K_S, w)$.
 - For all ind in $\text{DB}(w)$ in random order, initialize a counter $c \leftarrow 0$, then:
 - * Set $\text{xind} \leftarrow F_p(K_I, \text{ind})$, $z \leftarrow F_p(K_Z, w \parallel c)$ and $y \leftarrow \text{xind} \cdot z^{-1}$.
 - * Compute $\mathbf{e} \leftarrow \text{Enc}(K_e, \text{ind})$, and append (\mathbf{e}, y) to \mathbf{t} .
 - * Set $\text{xtag} \leftarrow g^{F_p(K_X, w) \cdot \text{xind}}$ and add xtag to XSet.
 - $\mathbf{T}[w] \leftarrow \mathbf{t}$.
- $(\text{TSet}, K_T) \leftarrow \text{TSetSetup}(\mathbf{T})$.
- Output the key $(K_S, K_X, K_I, K_Z, K_T)$ and $\text{EDB} = (\text{TSet}, \text{XSet})$.

Search protocol

- The client's input is the key $(K_S, K_X, K_I, K_Z, K_T)$ and query $\bar{w} = (w_1, \dots, w_n)$.
It sends to the server the message $(\text{stag}, \text{xtoken}[1], \text{xtoken}[2], \dots)$ defined as:
 - $\text{stag} \leftarrow \text{TSetGetTag}(K_T, w_1)$.
 - For $c = 1, 2 \dots$ and until server sends stop
 - * For $i = 2, \dots, n$, set $\text{xtoken}[c, i] \leftarrow g^{F_p(K_Z, w_1 \parallel c) \cdot F_p(K_X, w_i)}$
 - * Set $\text{xtoken}[c] = \text{xtoken}[c, 2], \dots, \text{xtoken}[c, n]$.
- The server has input $(\text{TSet}, \text{XSet})$. It responds as follows.
 - It sets $\mathbf{t} \leftarrow \text{TSetRetrieve}(\text{TSet}, \text{stag})$.
 - For $c = 1, \dots, |\mathbf{t}|$
 - * retrieve (\mathbf{e}, y) from the c -th tuple in \mathbf{t}
 - * if $\forall i = 2, \dots, n : \text{xtoken}[c, i]^y \in \text{XSet}$ then send \mathbf{e} to the client.
 - When last tuple in \mathbf{t} is reached, send stop to \mathcal{C} and halt.
- Client sets $K_e \leftarrow F(K_S, w_1)$; for each \mathbf{e} received, computes $\text{ind} \leftarrow \text{Dec}(K_e, \mathbf{e})$ and outputs ind .

Figure 3: OXT: OBLIVIOUS CROSS-TAGS PROTOCOL

These issues suggest a solution where we replace the function $f(\text{xtrap}, \cdot)$ (where $\text{xtrap} = F(K_X, w)$) with a form of *oblivious shared computation between client and server*. A first idea would be to use *blinded exponentiation* (as in Diffie-Hellman based oblivious PRF) in a group G of prime order p : Using a PRF F_p with range Z_p^* (and keys K_I, K_X), we derive a value $\text{xind} = F_p(K_I, \text{ind}) \in Z_p^*$ and define each xtag to be $g^{F_p(K_X, w) \cdot \text{xind}}$. The shared computation would proceed by the client first sending the value $g^{F_p(K_X, w_i) \cdot z}$ where $z \in Z_p^*$ is a blinding factor; the server would raise this to the power xind and finally the client would de-blind it by raising to the power $z^{-1} \bmod p$ to obtain $g^{F_p(K_X, w_i) \cdot \text{xind}}$. Unfortunately, this idea does not quite work as the server would learn $\text{xtag} = g^{F_p(K_X, w_i) \cdot \text{xind}}$ and from this, and its knowledge of xind , it would learn $g^{F_p(K_X, w_i)}$, which allows it to carry out an attack similar to the one against BXT. This also requires client-server interaction on a per-ind basis, a prohibitive cost.

In the design of OXT we address these two problems. The idea is to *precompute* (in EDBSetup) the blinding part of the oblivious computation and store it in the EDB. I.e., in each tuple corresponding to a keyword w and document xind , we store a blinded value $y_c = \text{xind} \cdot z_c^{-1}$, where z_c is an element in Z_p^* derived (via a PRF) from w and a tuple counter c (this counter, incremented for each tuple in \mathbf{t} associated with w , serves to ensure independent blinding values z).

During search, the server needs to compute the value $g^{F_p(K_X, w_i) \cdot \text{xind}}$ for each xind corresponding to a match for w_1 and then test these for membership in \mathbf{XSet} . To enable this, the client sends, for the c -th tuple in \mathbf{t} , a n -long array $\text{xtoken}[c]$ defined by $\text{xtoken}[c, i] := g^{F_p(K_X, w_i) \cdot z_c}$, $i = 1, \dots, n$, where z_c is the precomputed blinding derived by from w (via a PRF) and the tuple counter c . The server then performs the T-set search to get the results for w_1 , and filters the c -th result by testing if $\text{xtoken}[c, i]^{y_c} \in \mathbf{XSet}$ for all $i = 2, \dots, n$. This protocol is correct because

$$\text{xtoken}[c, i]^{y_c} = g^{F_p(K_X, w_i) \cdot z_c \cdot \text{xind} \cdot z_c^{-1}} = g^{F_p(K_X, w_i) \cdot \text{xind}},$$

meaning that the server correctly recomputes the pseudorandom values in the \mathbf{XSet} .

Putting these ideas together results in the OXT protocol of Figure 3. Note that the client sends the xtoken arrays (each holding several values of the form $g^{F_p(K_X, w_i) \cdot z}$) until instructed to stop by the server. There is no other communication from server to client (alternatively, server can send the number of elements in $\mathbf{TSet}(w)$ to the client who will respond with such number of xtoken arrays).⁷

Note that while the description above is intended to provide intuition for the protocol's design, assessing the security (leakage) of OXT is non-trivial, requiring an intricate security analysis that we provide in Section 5.

OXT consists of a single round of interaction, where the message sent by the client is of size proportional to $|\mathbf{DB}(w_1)|$,⁸ and the response to the client is minimal, consisting only of the result set (i.e., the set of encrypted ind's matching the query). The computational cost of OXT lies in the use of exponentiations, however, thanks to the use of very efficient elliptic curves (we only require the group to be DDH) and *fixed-base exponentiations*, this cost is practical even for very large databases as demonstrated by the performance numbers in Section 7.

OXT leaks much less information to the server than BXT. Indeed, since the server, call it \mathcal{S} , learns neither the ind values nor $\text{xtrap}_j, j = 2, \dots, n$, its ability to combine conjunctive terms from one query with terms from another query is significantly reduced. In particular, while in BXT \mathcal{S} learns the intersection between s-terms of any two queries, in OXT this is possible only in the

⁷ The same protocol supports single-keyword search (or 1-term conjunctions) by skipping the $c = 1, 2, \dots$ at both client and server, hence falling back to the SKS protocol of Figure 1.

⁸ For typical choices of w_1 , such message will be of small or moderate size. For large values of $|\mathbf{DB}(w_1)|$ one can cap the search to the first k tuples for a threshold k , say 1000. For example, in the case of a 3-term conjunction and xtag values of size 16 bytes, this will result in just 32 Kbyte message.

following case: the two queries can have different s-terms, but same x-term and there is a document containing both s-terms (the latter is possible since if the s-terms of two queries share a document `ind` and an x-term `xtrap` then the `xtag` value $f(\text{xtrap}, \text{ind})$ will be the same in both queries indicating that `ind` and `xtrap` are the same). The only other leakage via s-terms is that \mathcal{S} learns when two queries have the same s-term w_1 and the size of the set $\text{DB}(w_1)$. Finally, regrading intra-query leakage if \mathcal{C} responds with the values `xtagj`, $j = 2, \dots, n$, in the same order for all `ind`'s, then in case $n > 2$, \mathcal{S} learns the number of documents matching any sub-conjunction that includes w_1 and any subset of w_2, \dots, w_n . If, instead, \mathcal{C} randomly permutes the values `xtagj`, $j = 2, \dots, n$ before sending these values to \mathcal{S} , then \mathcal{S} learns the maximal number of satisfied terms per tuple in $\text{TSet}(w_1)$, but not the size of sets matching $w_1 \wedge w_i$, $i = 1, \dots, n$, or any other proper sub-conjunctions (except for what can be learned in conjunction with other leakage information). *In Section 5 we formally analyze the security of OXT making the above description of leakage precise.*

As noted before, even a leakage profile as the above that only reveals access patterns can still provide valuable information to an attacker that possesses prior information on the database and queries. We don't discuss here specific countermeasures for limiting the ability of an attacker to perform such statistical inference – see [23] for an example of potential masking techniques.

4 Processing Boolean Queries with OXT

We describe an extension to OXT that can handle *arbitrary* Boolean query expressions, but are efficient for a sub-class of expressions. We first consider conjunctions with negated terms.

Conjunction with negated terms. In the context of keyword-based search we refer to a “negated term” as one that returns documents that *do not contain* the given keyword. Consider a conjunction of n terms in which there is *at least one non-negated term*. To search for such a conjunction we modify OXT as follows.

- The Client chooses one of the non-negated terms as the s-term, and computes `stag` and `xtraps` (i.e. the `xtoken` arrays) as in OXT. It then sends the `stag` and the `xtraps` to the Server, but indicating which `xtraps` are for negated terms.
- The Server's computation is exactly as in OXT, except that for negated terms it checks if the corresponding $(\text{xtoken}[c, i])^y$ is *not* in the `XSet` (i.e. instead of checking if it is in the `XSet`).

A conjunction where *all* terms are negated cannot be executed as above since the Client cannot choose a negated term as s-term. We will see how we can accommodate, though inefficiently, such searches below. (Note that in most cases a negated term will have a very large number of matching documents; if this is the case for all the terms in the conjunction then efficient search, even on plaintext data, is not possible.)

Boolean expressions in Searchable Normal Form (SNF). We say that a Boolean expression in n terms is in *Searchable Normal Form* (SNF) if it is of the form $w_1 \wedge \phi(w_2, \dots, w_n)$ where ϕ is an *arbitrary* Boolean formula (e.g., “ $w_1 \wedge (w_2 \vee w_3 \vee \neg w_4)$ ”). Protocol OXT can be extended to answer such queries; the needed modifications are similar to those described above for the case of conjunctions with at least one non-negated term (a special case of SNF). Specifically, on input a query of the form $w_1 \wedge \phi(w_2, \dots, w_n)$, the client creates a modified boolean expression $\hat{\phi}$ in new boolean variables v_i ($i = 2, \dots, n$), which is just ϕ but with each w_i replaced by v_i . Thus, the client uses w_1 as the s-term and computes its `stag` as in OXT, and computes the `xtrap` (i.e. the `xtoken` array) for all the other terms w_i ($i > 1$). It then sends the `stag` and the `xtraps` in the order of their index. It also sends the Server the above modified boolean expression $\hat{\phi}$.

The Server fetches the TSet corresponding to the stag as in OXT. It also computes the xtag corresponding to each x-term, also as in OXT. But, it decides on sending (to the Client) the encrypted ind corresponding to each tuple in the TSet based on the following computation (which is the only different part from OXT): for each $i = 2, \dots, n$, the Server treats the variable v_i as a boolean variable and sets it to the truth value of the expression $(\text{xtoken}[c, i])^y \in \text{XSet}$. Then it evaluates the expression $\phi(v_2, \dots, v_n)$. If the result is true, it returns the e value in that tuple to the Client.

Responding to arbitrary Boolean queries. OXT can be also be extended to answer *any* Boolean query by adding to the database a field TRUE which all documents satisfy. Then a search for any expression $\phi(w_1, \dots, w_n)$ can be implemented as “TRUE $\wedge \phi(w_1, \dots, w_n)$ ”, which is in SNF and can be searched as in the SNF case above. Clearly, this will take time linear in the number of documents but it can be implemented if such functionality is considered worth the search complexity.

Disjunctions. A disjunction of expressions that have efficient search under OXT or the extended OXT as described above can also be searched efficiently by running OXT (or the extended OXT) in each disjunct separately, and the client requesting the union of the sets of decrypted ind’s for each disjunct.

5 Security Analysis of OXT

In this section we describe the OXT leakage profile \mathcal{L} and analyze its security. While our ultimate goal is to prove adaptive security for boolean queries, we start by analyzing *non-adaptive* security for the special case where all queries are conjunctions of two keywords. This special case already captures essentially all of the difficulties in the analysis and dispenses with several distracting notational issues. Then we show how to extend the leakage profile to boolean queries (in particular for conjunctions with any number of terms) in Section 5.3, where we also analyze adaptive security.

5.1 Hardness assumptions

Decision Diffie-Hellman. Let $G = G_\lambda$ be a prime order cyclic group of order $p = p(\lambda)$ generated by g . We say that the *decision Diffie-Hellman (DDH) assumption* holds in G if $\text{Adv}_{G,A}^{\text{ddh}}(\lambda)$ is negligible for all efficient adversaries A , where

$$\text{Adv}_{G,A}^{\text{ddh}}(\lambda) = \Pr[A(g, g^a, g^b, g^{ab}) = 1] - \Pr[A(g, g^a, g^b, g^c) = 1]$$

where the probability is over the randomness of A and uniformly chosen a, b, c from Z_p^* .

For vectors $\mathbf{a} \in (Z_p^*)^\alpha, \mathbf{b} \in (Z_p^*)^\beta$ let $g^{\mathbf{a}} = (g^{a[1]}, \dots, g^{a[\alpha]}) \in G^\alpha$ and $g^{\mathbf{a}^T \mathbf{b}}$ be the matrix in $G^{\alpha \times \beta}$ where the (i, j) -th entry is $g^{a[i] \cdot b[j]}$. We will use the following standard lemma in our security proof.

Lemma 3 *Suppose the DDH assumption holds for in G . Then, for any integers α, β (polynomial in λ) any efficient adversary A , we have*

$$\Pr[A(g, g^{\mathbf{a}}, g^{\mathbf{b}}, g^{\mathbf{a}^T \mathbf{b}}) = 1] - \Pr[A(g, g^{\mathbf{a}}, g^{\mathbf{b}}, \mathbf{M}) = 1] \leq \text{neg}(\lambda),$$

where \mathbf{a} is uniform over $(Z_p^*)^\alpha$, \mathbf{b} is uniform over $(Z_p^*)^\beta$, and \mathbf{M} is uniform over $G^{\alpha \times \beta}$.

PRF Security. Let X and Y be sets, and let $F : \{0, 1\}^\lambda \times X \rightarrow Y$ be a function. We say that F is a *pseudorandom function (PRF)* if for all efficient adversaries A , $\text{Adv}_{F,A}^{\text{prf}}(\lambda)$ is negligible, where

$$\text{Adv}_{F,A}^{\text{prf}}(\lambda) = \Pr[A^{F(K, \cdot)}(1^\lambda) = 1] - \Pr[A^{f(\cdot)}(1^k) = 1]$$

where the probability is over the randomness of A , $K \xleftarrow{\$} \{0, 1\}^\lambda$, and $f \xleftarrow{\$} \text{Fun}(X, Y)$.

IND-CPA Encryption Security. A symmetric encryption scheme (Enc, Dec) is a pair of algorithms, the first randomized and the second deterministic. Enc takes as input a key $K \in \{0, 1\}^\lambda$ and a message $M \in \{0, 1\}^*$, and outputs a ciphertext C . Dec takes as input a key $K \in \{0, 1\}^\lambda$ and a ciphertext C and outputs a message M . We require the usual correctness for all possible keys K and messages M .

We say that $\Sigma = (\text{Enc}, \text{Dec})$ is *IND-CPA secure* if for all efficient adversaries A $\text{Adv}_{\Sigma, A}^{\text{ind-cpa}}(\lambda)$ is negligible, where

$$\text{Adv}_{\Sigma, A}^{\text{ind-cpa}}(\lambda) = \Pr[A^{\mathcal{O}(K, 0, \cdot, \cdot)}(1^\lambda) = 1] - \Pr[A^{\mathcal{O}(K, 1, \cdot, \cdot)}(1^\lambda) = 1],$$

where K is chosen at random from $\{0, 1\}^\lambda$ and the oracle $\mathcal{O}(K, b, M_0, M_1)$ returns \perp if $|M_0| \neq |M_1|$, and otherwise it samples $\text{Enc}(K, M_b)$ and returns the result.

5.2 Warm up: Analysis for non-adaptive 2-conjunctions

We start by describing the function \mathcal{L}_{ext} that describes the leakage of the OXT (beyond the T-set leakage) protocol under a simpler scenario where all queries are non-adaptive and are for conjunctions of two keywords only. This setting already confronts most of the intuitive difficulties in our leakage profile and proof, but is less cluttered than the full setting, which we present in the next subsection.

Below our security theorem will show that this, in addition to the leakage from \mathcal{L}_T which is defined by the T-set implementation, is *all* of the information leaked by our protocol.

We represent a sequence of Q non-adaptive 2-conjunction queries by $\mathbf{q} = (\mathbf{s}, \mathbf{x})$ where an individual query is a 2-term conjunction $\mathbf{s}[i] \wedge \mathbf{x}[i]$ which we write as $\mathbf{q}[i] = (\mathbf{s}[i], \mathbf{x}[i])$. $\mathcal{L}_{\text{ext}}(\text{DB}, \mathbf{q})$ gets $\text{DB} = (\text{ind}_i, W_i)_{i=1}^d$ and $\mathbf{q} = (\mathbf{s}, \mathbf{x})$ as input and outputs $(N, \bar{\mathbf{s}}, \text{SP}, \text{RP}, \text{IP})$, which are defined below.

- $N = \sum_{i=1}^d |W_i|$ is the total number of appearances of keywords in documents.
- $\bar{\mathbf{s}} \in [m]^Q$ is the *equality pattern* of $\mathbf{s} \in W^Q$ indicating which queries have the equal s-terms. Formally, $\bar{\mathbf{s}} \in [m]^Q$ is formed by assigning each keyword an integer in $[m]$ determined by the order of appearance in \mathbf{s} . For example, if $\mathbf{s} = (a, a, b, c, a, c)$ then $\bar{\mathbf{s}} = (1, 1, 2, 3, 1, 3)$. To compute $\bar{\mathbf{s}}[i]$ one finds the least j such that $\mathbf{s}[j] = \mathbf{s}[i]$ and then lets $\bar{\mathbf{s}}[i] = |\{\mathbf{s}[1], \dots, \mathbf{s}[j]\}|$ be the number of unique keywords appearing at indices less than or equal to j .
- SP is the *size pattern* of the queries, which is the number of documents matching the first keyword in each query. Formally, $\text{SP} \in [d]^Q$ and $\text{SP}[i] = |\text{DB}(\mathbf{s}[i])|$.
- RP is the *results pattern* of the queries, which are the indices of documents matching the entire conjunction. Formally, RP is vector of size Q with $\text{RP}[i] = \text{DB}(\mathbf{s}[i]) \cap \text{DB}(\mathbf{x}[i])$ for each i .
- IP is the *conditional intersection pattern*, which is formally a Q by Q table defined by

$$\text{IP}[i, j] = \begin{cases} \text{DB}(\mathbf{s}[i]) \cap \text{DB}(\mathbf{s}[j]) & \text{if } i \neq j \text{ and } \mathbf{x}[i] = \mathbf{x}[j] \\ \emptyset & \text{otherwise} \end{cases}$$

Understanding the leakage components. The parameter N can be replaced with an upper bound given by the total size of EDB but leaking such a bound is unavoidable. The equality pattern

\bar{s} leaks repetitions in the s-term of different queries; this is a consequence of our optimized search that singles out the s-term in the query. This leakage can be mitigated by having more than one TSet per keyword and the client choosing different incarnations of the Tset for queries with repeated s-terms. SP leaks the number of documents satisfying the s-term in a query and is also a direct consequence of our approach of optimizing search time via s-terms; it can be mitigated by providing an upper bound on the number of documents rather than an exact count by artificially increasing Tset sizes. RP is the result of the query and therefore no real leakage in the context of SSE. Finally, the IP component is the most subtle and it means that if two queries have different s-terms but same x-term, then *if there is a document satisfying both s-terms* then the set of indexes matching both s-terms is leaked (if no document matching both s-terms exist then nothing is leaked). It can be seen as the price we pay for the rich functionality enabled by our x-terms and XSet approach that allows for the computation of arbitrary boolean queries. Note, however, that since the s-terms are meant to be the least-frequently matched keywords, the number of instances with s-terms having a non-empty intersection of documents containing both terms can be minimized. Moreover, in searches where the s-term is a unique per-document term (e.g., a last name on a database with a last-name field) the IP leaks nothing about the searches or database, as all the S_i are empty.

It also helps to compare the above leakage to the leakage incurred by the naïve solution that simply sends tokens for each keyword in the conjunction. In that case, the server would learn $DB(s[i])$ and $DB(x[i])$ for every i , as it would see the results of the individual searches (instead of just the results pattern, which would be ideal). This means it would learn the relationships between all of the documents matching the individual keywords, and this is information is exactly what OXT is designed to minimize.

We also remark that d , the number of documents in DB is not leaked in the above profile but would like be leaked in an implementation that stores the encrypted documents. This is outside the formal model analyzed here but could be easily incorporated at the cost of extra notation, as some prior works did.

OXT security theorem for two-term conjunctions. Theorem 5 below that states the security of OXT for two-term conjunctions requires the following lemma that establishes a bound on the correctness of the protocol, which for OXT means the probability that a document that does not match a query will be returned to the client. The next lemma shows that this happens only with negligible probability. See Section 5.1 for standard assumptions (PRF, DDH) used in the theorem.

Lemma 4 *For every adversary A there exists adversaries B and B' which run in essentially the same time as A , such that*

$$\Pr[\text{Cor}_A^{\text{OXT}}(\lambda) = 1] \leq 2 \cdot \text{Adv}_{F_p, B}^{\text{prf}}(\lambda) + \text{AdvCor}_{B'}^{\Pi}(\lambda) + N^2/(p-1) + N/p,$$

where $N = \sum_{i=1}^d |W_i|$ is the total number of appearances of keywords in all documents, p is the order of the group G , and Π is the T-set implementation.

The proof appears in Appendix A. Intuitively, the only way a failure in correctness can occur is if either the T-set fails or if two keyword-identifier pairs from DB map to the same xtag using $F_p(K_I, \cdot)$ and $F_p(K_X, \cdot)$. But since these values are pseudorandom over a large set, the chance of this happening is negligible.

The statement of our generic security theorem for a generic T-set is somewhat cumbersome due to an issue with how leakages compose, and is given in Section B. For the case of OXT that uses our T-set implementation from Section 6, which is of primary interest, we have that \mathcal{L}_{oxt} constitutes all of the leakage.

Theorem 5 *Let \mathcal{L}_{OXT} be as defined above, and suppose that the T-set implementation Σ from Section 6. Then SSE scheme OXT is \mathcal{L}_{OXT} -semantically-secure against non-adaptive attacks where all queries are 2-conjunctions, assuming that the DDH assumption holds in G , that F and F_p are secure PRFs, that (Enc, Dec) is an IND-CPA secure symmetric encryption scheme, and the conditions from Theorem 7 hold.*

Proof sketch. The proof of the theorem is delicate and lengthy, and is given in Appendix B for the more general formulation in Theorem 9. To get some intuition for why the scheme is secure, we start by examining why each of the outputs of \mathcal{L} is necessary for a correct simulation. Of course, this does nothing to show that they are *sufficient* for simulation, but it will be easier to see why this is all of the leakage once their purpose is motivated.

The size of the XSet is equal to the value N leaked. The equality pattern for \mathbf{s} , $\bar{\mathbf{s}}$, (or something computationally equivalent to it) is necessary due to the fact that the **stag** values are deterministic, so a server can observe repetitions of **stag** values to determine if $\mathbf{s}[i] = \mathbf{s}[j]$ for all i, j . The size pattern is also necessary as the server will always learn the number of matches for the first keyword in the conjunction by observing the number of tuples returned by the T-set. We include the results pattern to enable the simulator to produce the client results for queries in way consistent the conditional intersection pattern.

The final and most subtle part of the leakage is the conditional intersection pattern IP. The IP is present in the leakage because of the following passive attack. During the computation of the search protocol, the values tested for membership in the XSet by the server have the form $g^{F_p(K_X, w_i) \cdot F_p(K_I, \text{ind})}$, where w_i is the i -th keyword from a search and ind is an identifier for a document that matched the s-term (but may or may not match the remaining keywords). The leakage comes from the fact that the values will sometimes *repeat* (when two queries repeat the same w_i and match the same ind with their s-terms) and they are all known to the adversary. The IP describes exactly what can be learned from this attack.

Our proof makes formal the claim that the output of \mathcal{L} is sufficient for a simulation. We outline a few of the technical hurdles in the proof without dealing with the details here. For this discussion, we assume that reductions to PRF security and encryption security go through easily, allowing us to treat PRF outputs as random and un-opened ciphertexts as encryptions of zeros.

We first handle the information leaked by the XSet. An unbounded adversary could compute the discrete logarithms of the XSet elements and derive information about which documents match which keywords. We want to show however that a poly-time adversary learns nothing from the XSet due to the assumed hardness of the DDH problem. Formally, we need to show that we can replace the elements of XSet with random elements that carry no information about the database, but there is a technical difficulty: some of the exponents (specifically, the xind values) that will play the roll of hidden exponents in the DDH reduction are used in the computation of the xtrap values, and these are revealed in the transcripts. A careful rearrangement of the game computation will show that this is not as bad as it seems, because the xind values are “blinded out” by the z values. We stress that this requires some care, because the z values are also used twice, and we circumvent this circularity by computing the XSet first and then computing the transcripts “backwards” in way that is consistent with the XSet. Now a reduction to DDH becomes clear, as the XSet values can be dropped in obviously as real-or-random group elements.

With the XSet leakage eliminated, the rest of the work is in showing that the simulator can arrange for a correct-looking pattern of “repeats” in the documents matched and in the values tested against the XSet. While riddled with details, this is intuitively a rather straightforward task that is carried out in the latter games of the proof.

5.3 Analysis for boolean queries

In this section we give the leakage profile for OXT under adaptive attacks where any boolean query in searchable normal form (SNF), i.e., of the form $\psi(s, x_1, \dots, x_n) = s \wedge \phi(x_1, \dots, x_n)$, is allowed (n can vary between queries). In particular, this leakage profile extends the one in the previous section to conjunctions with any number of terms. As before, for simplicity we describe the leakage profile when OXT is used with our specific T-set implementation from Section 6.

The new leakage profile \mathcal{L} is a stateful algorithm that responds to inputs as defined in Definition 1. On the initial input DB, it outputs $N = \sum_{i=1}^d |W_i|$. Later inputs consist of a vector of queries $\mathbf{q} = (\Phi, \mathbf{s}, \mathbf{x}_1, \dots, \mathbf{x}_n)$, where Φ is a vector of boolean formulae and $\mathbf{s}, \mathbf{x}_1, \dots, \mathbf{x}_n$ are vectors of keywords, all populated according to the queries issued so far in the straightforward way⁹. The leakage function \mathcal{L} outputs $(N, \Phi, \bar{\mathbf{s}}, \mathbf{SP}, \mathbf{XP}, \mathbf{RP}, \mathbf{IP})$, where $N, \bar{\mathbf{s}}, \mathbf{SP}$ are computed as in Section 5.2 and Φ is from \mathbf{q} . It computes the remaining elements as follows:

- The vector \mathbf{XP} has $\mathbf{XP}[i]$ set to the number of x-terms in the i -th query.
- The *results pattern* \mathbf{RP} is a vector, where for each i , $\mathbf{RP}[i]$ is the identifiers matching the i -th query.
- The *conditional intersection pattern* \mathbf{IP} is intuitively the \mathbf{IP} from the 2-conjunction case for every pair of x-term indices. Formally, \mathbf{IP} is indexed by (i, j, α, β) where $1 \leq i, j \leq Q$ and $1 \leq \alpha, \beta \leq n$ (where n is the maximum number of x-terms in any query). The entries in \mathbf{IP} are defined by

$$\mathbf{IP}[i, j, \alpha, \beta] = \begin{cases} \text{DB}(\mathbf{s}[i]) \cap \text{DB}(\mathbf{s}[j]) & \text{if } i \neq j, \alpha \neq \beta \text{ and } \mathbf{x}_\alpha[i] = \mathbf{x}_\beta[j] \\ \emptyset & \text{otherwise} \end{cases}$$

Theorem 6 *Let \mathcal{L} be the leakage function defined above. Then OXT, when implemented using the T-set implemenation from Section 6, is \mathcal{L} -semantically-secure against adaptive attacks, assuming that the DDH assumption holds in G , that F and F_p are secure PRFs, that (Enc, Dec) is an IND-CPA secure symmetric encryption scheme, and that the conditions in Theorem 7 are satisfied.*

The proof of this theorem is an extension of the non-adaptive theorem in two ways. First, the proof must handle the simulation of more complicated general queries, and second it must do this adaptively. Handling general queries introduces few complications because the leakage of Φ and the new \mathbf{IP} gives the simulator enough information to properly program the results and the \mathbf{XSet} in a way that generalizes what is done in the proof for 2-conjunctions. Thus we focus on how to handle adaptivity.

The adaptive simulator works as follows. To generate $\text{EDB} = (\text{TSet}, \text{XSet})$, the simulator invokes the adaptive TSet simulator from the proof of Theorem 7 on input N and it generates XSet by choosing N random group elements and adding them to XSet .

To simulate queries responses, the simulator will adaptively “assign” elements of the XSet to keyword-ind pairs. This is in contrast with the non-adaptive simulator, which achieved this by initializing the array H and then adding the elements to the XSet as determined by the leakage. Here, the simulator is choosing the XSet values, and then initializing H entries adaptively, using either elements of XSet (with the proper repetitions) or independent random elements (again, with the proper repetitions).

⁹Formally, n is the maximum number of x-terms in any query, with smaller queries padded up with a special dummy symbol. Hence, if the i -th query uses boolean formula ϕ on $n' + 1$ inputs then we have $\Phi[i] = \phi$ and the input is represented by $(\mathbf{s}[i], \mathbf{x}_1[i], \dots, \mathbf{x}_n[i])$ where the last $n - n'$ terms are set to the dummy symbol.

6 T-Set Instantiation

We show an efficient instantiation of a T-set, whose syntax and security are defined in Section 2. We instantiate a T-set as a hash table with B buckets of size S each. The $\text{TSetSetup}(\mathbf{T})$ procedure sets the parameters B and S depending on the total number $N = \sum_{w \in \mathbf{W}} |\mathbf{T}[w]|$ of tuples in \mathbf{T} in such a way so that (1) the probability of an overflow of any bucket after storing N elements in this hash table is a sufficiently small constant; and (2) the total size $B \cdot S$ of the hash table is $O(N)$.

Figure 4 shows our T-set instantiation $\Sigma = (\text{TSetSetup}, \text{TSetGetTag}, \text{TSetRetrieve})$. We use the following notation: λ is a security parameter, $n(\lambda)$ is the bit length of strings s_i in each vector \mathbf{T} , **record** is a record type with two fields, **label**, storing bit strings of length λ , and **value**, storing bit strings of length $n(\lambda) + 1$. Each entry $\text{TSet}[i]$, for $i = 1, \dots, B$, of array TSet , is itself an array of S records of type **record**. We denote the j -th element in the array $\text{TSet}[i]$ as $\text{TSet}[i, j]$. The protocol uses PRFs F, \bar{F} and a hash function H , which for adaptive security is modeled as a random oracle in the analysis, where F maps integers between 1 and $\max_w |D[w]|$ onto $\{0, 1\}^\lambda$, \bar{F} maps the set W onto the key space of F , and H maps $\{0, 1\}^\lambda$ onto $\{1, \dots, B\} \times \{0, 1\}^\lambda \times \{0, 1\}^{n(\lambda)+1}$.

T-Set correctness. Note that procedure TSetSetup places the i -th element s_i of $\mathbf{T}[w]$ into bucket $\text{TSet}[b]$ where b is the first element in tuple $(b, L, K) = H(F(\text{stag}, i))$, for $\text{stag} = \bar{F}(K_T, w)$. Therefore the search procedure $\text{TSetRetrieve}(\text{TSet}, \text{stag})$ will recover the same bucket $\mathbf{B} = \text{TSet}[b]$ for $\text{stag} \leftarrow \text{TSetGetTag}(K_T, w)$. Since that i -th element s_i was placed in a record r in $\text{TSet}[b]$ s.t. $r.\text{label} = L$, the same record is found in the TSetRetrieve procedure as $\mathbf{B}[j]$, unless bucket \mathbf{B} contains some other record with the same label field L . However, the probability that this ever happens can be bounded by $BS^22^{-\lambda}$ (plus at most a negligible quantity bounding the maximal distinguishing advantage between pseudorandom functions F and \bar{F} and true random functions). If the correct record r is found then the TSetRetrieve procedure uncovers the correct bit β and string s_i from $r.\text{value}$ because it is xored by the same one-time pad K , the third element in the same triple (b, L, K) , with which $(\beta|s_i)$ was encrypted in $r.\text{value}$ during TSetSetup . Since $\beta = 0$ if and only if s_i is the last string in list $\mathbf{T}[w]$, TSetRetrieve finds all elements of $\mathbf{T}[w]$ and terminates when $\mathbf{t} = \mathbf{T}[w]$.

Note that since a single attempt to build the TSet table by procedure TSetSetup takes $O(m)$ steps, and parameters B, S are set so that a single attempt fails with at most a constant probability over the choice of the key K_T of PRF \bar{F} , the expected time for the TSetSetup procedure is $O(m)$.

T-Set security. The only leakage incurred by the T-set instantiation of Figure 4 is the total number $N = \sum_{w \in \mathbf{W}} |\mathbf{T}[w]|$ of tuples in \mathbf{T} , which in our SSE applications equals the total number $\sum_{w \in \mathbf{W}} |\text{DB}(w)|$ of keyword occurrences in database DB.

Theorem 7 *For any keyword sequence \mathbf{q} , including an empty sequence, define $\mathcal{L}_T(\mathbf{T}, \mathbf{q})$ as $\sum_{w \in \mathbf{W}} |\mathbf{T}(w)|$, where \mathbf{W} is the set of keywords used by \mathbf{T} . The T-set instantiation Σ in Figure 4 is \mathcal{L}_T -adaptively-secure assuming that F and \bar{F} are secure PRFs and that H is a random oracle.*

Dispensing of the random oracle. Theorem 7 is the only result in this paper that uses the random oracle model. The model is used to prove adaptive security (and indeed to avoid the lower bound of [14] on token lengths). For non-adaptive security, we can prove security for a version of the protocol with H replaced with another PRF (and an additional key stored at the client). Moreover, it is possible to prove adaptive security without a random oracle at the cost of added communication. For this we dispense of H and let the client send the values $F(\text{stag}, i)$ instead of stag . This increases the amount of communication, but when used with our OXT protocol it is not a dramatic increase. As with OXT, the server can send a “stop” message when the client has sent enough $F(\text{stag}, i)$ values.

<p><u>TSetSetup(T)</u></p> <ul style="list-style-type: none"> • Initialize an array TSet of size B whose every element is an array of S records of type record. • Initialize an array Free of size B whose elements are integer sets, initially all set to $\{1, \dots, S\}$. • Choose a random key K_T of PRF \bar{F}. • Let W be the set of keywords in DB. For every $w \in W$ do the following: <ul style="list-style-type: none"> – Set $\text{stag} \leftarrow \bar{F}(K_T, w)$ and $\mathbf{t} \leftarrow \mathbf{T}[w]$. – For each $i = 1, \dots, \mathbf{t}$, set s_i as the i-th string in \mathbf{t}, and perform the following steps: <ul style="list-style-type: none"> * Set $(b, L, K) \leftarrow H(F(\text{stag}, i))$. * If Free[b] is an empty set, restart TSetSetup(T) with fresh key K_T. * Choose $j \xleftarrow{\\$} \text{Free}[b]$ and remove j from set Free[b], i.e. set $\text{Free}[b] \leftarrow \text{Free}[b] \setminus \{j\}$. * Set bit β as 1 if $i < \mathbf{t}$ and 0 if $i = \mathbf{t}$. * Set $\text{TSet}[b, j].\text{label} \leftarrow L$ and $\text{TSet}[b, j].\text{value} \leftarrow (\beta s_i) \oplus K$. • Output (TSet, K_T). <p><u>TSetGetTag(K_T, w)</u></p> <ul style="list-style-type: none"> • Output $\text{stag} \leftarrow \bar{F}(K_T, w)$ <p><u>TSetRetrieve(TSet, stag)</u></p> <ul style="list-style-type: none"> • Initialize \mathbf{t} as an empty list, bit β as 1, and counter i as 1. • Repeat the following loop while $\beta = 1$: <ul style="list-style-type: none"> – Set $(b, L, K) \leftarrow H(F(\text{stag}, i))$ and retrieve an array $\mathbf{B} \leftarrow \text{TSet}[b]$ – Search for index $j \in \{1, \dots, S\}$ s.t. $\mathbf{B}[j].\text{label} = L$. – Let $v \leftarrow \mathbf{B}[j].\text{value} \oplus K$. Let β be the first bit of v, and s the remaining $n(\lambda)$ bits of v. – Add string s to the list \mathbf{t} and increment i. • Output \mathbf{t}.

Figure 4: T-SET INSTANTIATION Σ

Achieving optimal leakage. The leakage $\mathcal{L}_T(\mathbf{T}, \mathbf{q}) = \sum_{w \in W} |\mathbf{T}(w)|$ can be relaxed to only disclose an upper bound on the latter quantity by filling in all empty locations in all buckets with (pseudo) random values. In that way, the only leakage is the size of the full hash table which is optimal.

Proof sketch. To prove Theorem 7, we start with the simulator algorithm S . The simulator's initial input is $N = \mathcal{L}_T(\mathbf{T}) = \sum_{w \in W} |\text{DB}(w)|$, and for each query q the adversary makes, the simulator gets the corresponding vector $\mathbf{T}[q]$, and no additional leakage. At initialization $S(N)$ creates **TSet** as a $B \times S$ table (note that parameters B, S are determined by N) just like the TSetSetup algorithm, but S populates this table by choosing m random locations in it and filling them up with random entries. In other words, for each $i = 1, \dots, m$, S randomly chooses a block

b and a free location j within this block, and assigns $\text{TSet}[b, j]$ to a record whose label and value fields are random bitstrings of appropriate length. S hands this table TSet over to A , but locally it marks all these m records in TSet as “unrevealed”. In the Random Oracle Model for hash function H , S also simulates the outputs of H for A , initially setting $H(x)$ for every query x at random. On every query q which adversary A makes, S receives a list $\mathbf{t} = (s_1, \dots, s_{|\mathbf{t}|})$ s.t. $\mathbf{t} = \mathbf{T}[q]$. (Note that function \mathcal{L}_T does not include any further leakage to S .) $S(\mathbf{t})$ picks value stag' at random in the key space of PRF F , chooses $|\mathbf{t}|$ random unrevealed records $r_1, \dots, r_{|\mathbf{t}|}$ in TSet , locally marking them all as revealed, and for each $i = 1, \dots, |\mathbf{t}|$, S sets $H(F(\text{stag}', i))$ to a tuple (b, L, K) defined as follows: Bucket b is set to the bucket where record r_i is, label L is set to $r_i.\text{label}$, and key K is set to $(r_i.\text{value}) \oplus (\beta |s_i|)$ where β is set to 1 if $i < |\mathbf{t}|$ and 0 if $i = |\mathbf{t}|$. If H was queried on any of the points $F(\text{stag}', i)$ before, S aborts. Otherwise, S outputs stag' to A as its response to $\mathbf{t} = \mathbf{T}[q]$.

First, note that by the PRF property of \bar{F} , the real game is indistinguishable from a modification in which values stag'_w for all $w \in \mathcal{W}$ are replaced by random elements of the range of \bar{F} , which is the keyspace of PRF F . Secondly, because of the PRF property of F , this keyspace must be large enough so that there only negligible probability that A queries H on (stag'_w, i) for any $w \in \mathcal{W}$ and any i before A sees the corresponding stag'_w value. Finally, since records corresponding to any given $\mathbf{T}[q]$ sequence are assigned at random to the TSet table, under the condition that H isn't queried on any point (stag', i) before S releases the corresponding stag' value to A , the adversary's view of each $(b, L, K) = H(F(\text{stag}', i))$ tuple and the corresponding record r_i is identical in the above simulation and the execution modified above, because in both cases b, L, K are random, $r_i.\text{label} = L$ and $r_i.\text{value} = K \oplus (\beta, s_i)$ where $\beta = 1$ for $i < |\mathbf{t}|$ and 0 for $i = |\mathbf{t}|$. This completes the argument for adaptive security of our T-set instantiation.

Overflow probability. Here we provide some examples of overflow probabilities for a given number of buckets and their size. Consider a hash table with B buckets, each with space for S equal-sized items. Let there be N items, which are inserted randomly and independently into the hash table buckets. Let X_i^j denote the indicator variable for the event that the i -th item was inserted in j -th bucket. Then, for each bucket j , the event W^j that it overflows is given by $\sum_{i=1..N} X_i^j > S$. First note that for all i, j , $\Pr[X_i^j = 1] = 1/B$. Thus, by linearity of expectation, for each j , $\mathbf{E}[\sum_{i=1..N} X_i^j] = N/B$. If k , the *space-overhead*, is defined as $B \cdot S/N$, then the above expectation is S/k . Now by Chernoff bound [2], for $k \geq 1$,

$$\Pr[W^j] = \Pr\left[\sum_{i=1..N} X_i^j > k \cdot (S/k)\right] < [e^{k-1} k^{-k}]^{S/k} < (e^{1-1/k}/k)^S.$$

By the union bound, the probability that any of the B buckets overflows is at most $B \cdot (e/k)^S$ which in turn equals $\frac{Nk}{S} \cdot (e^{1-1/k}/k)^S$. To illustrate the overflow probability with specific numerical examples we consider $N = 2^{30}$. Then with $k = 3, S = 80$ one gets overflow probability of less than 2^{-21} (and with $S = 100$ less than 2^{-37}). With $k = 2, S = 160$ the probability is 2^{-21} (and with $S = 200$ less than 2^{-32}). With $k = 1.5, S = 400$ the probability is about 2^{-20} (and with $S = 500$ less than 2^{-30}). Finally, if one uses moderately large buckets, then even with a mere 10% increase in database size one can get very low overflow probabilities; e.g., with $k = 1.1, S = 6000$ the probability is 2^{-20} (and with $S = 8000$ less than 2^{-33}).

Note that this is the probability that pre-processing fails and needs to be re-started; this probability is not adversarially controlled and has no security consequences.

7 OXT Implementation and Experimental Results

This section reports the status of the OXT implementation and some results from using our prototype with three different data sets. The latency and scalability measurements included here should be viewed as providing empirical proof to the performance and scalability claims made earlier in the paper. The section starts with a brief description of the prototype that focuses on the TSet and XSet generation and their use in the query processing, which are the biggest challenges in implementing OXT. Next, the data sets are described followed by latency and scalability measurements. The section concludes with our plans for improving the protocol performance.

Prototype. The prototype consists of two main components, one for data set pre-processing and the other for query execution, in addition to a relatively simple client program used for testing and performance measurements. The three programs use the same cryptographic primitives, which leverage the OpenSSL 1.0.1 library. To realize the cryptographic primitives mentioned in Section 2, we use the NIST 224p elliptic curves for group operations, AES-FFX [5] for permutations and, depending on platform characteristics, AES-CMAC or HMAC for PRFs. Encryption optimized for small storage overhead is achieved with a stateful and position-based variant of AES in CTR mode. The overall C code, including test programs, measures roughly 16k lines.

The pre-processing generates the TSet and XSet data structures (or EDB) for the query component as flat files together with the meta-data describing the EDB. During initialization, the query-execution component reads the EDB meta-data, retrieves the RAM-resident structures from the flat files and, subsequently, answers queries from clients, one query at a time.

The design of the TSet and XSet data structures is driven by the scalability requirement, to data sets well beyond to RAM capacity of high-end servers or even of server clusters, and by our goal of minimizing query execution time. To achieve the latter, the design maximizes the use of the available RAM and leverages domain characteristics, such as tolerance to small false positive rates. In addition, the data structures had to be designed to properly handle the protocol-induced, strictly random disk I/O operations without introducing implementation-induced leakage.

For scalability, the TSet is realized as a disk-resident paged hash table. To maximize data retrieval per disk I/O operation, each tuple list $T[w]$ in the TSet is segmented into fixed-size blocks of tuples keyed by a tag stag_c . This tag is derived by a PRF from the list's stag and a segment counter. Tuple blocks are grouped, based on their tags, in (hash table) pages, which occupy contiguous disk areas. The unused blocks in each page are filled with random bits to make them indistinguishable from the blocks used to store tuples. The used/unused type of a block remains unknown to the query execution component until the block's tag gets generated in a client query.

The page size is set to the stride size of the underlying RAID-5 array. Similarly, the (intermediate) ext4 file system is configured based on the RAID characteristics. The query execution component uses page-level direct I/O to prevent buffer cache pollution in the OS, as the hash table pages are inherently uncachable. In addition, this component parallelizes disk accesses using asynchronous I/O (`aio_*` system calls). As a result, each page access is handled by a single hard disk drive and each drive has several outstanding page requests, which achieves the maximum possible parallelism and I/O depth in the storage system. Note that due to the inherent (and intended) random nature of the TSet accesses, neither the underlying file system nor the storage subsystem can leverage standard acceleration techniques, such as pre-fetching, to improve system performance.

The much smaller XSet is realized as a RAM-resident Bloom filter [6], which enables the sizing the false positive rate to the lowest value that the server's RAM allows. In the the current prototype, we set the false positive rate to 2^{-20} and the resulting Bloom filter XSet still occupies only a small fraction of our server RAM. During query execution, disk accesses are overlapped with the xtag generation and Bloom filter query operations, which are both parallelized across all the available

CPU cores.

Data Sets. To show the practical viability of our solution we run tests on three data sets: a 100,000 record (attribute-value) database synthesized from census data, where each record is treated as a separate document d , attribute values are atomic, and where each pair $(attribute_i, value)$ is treated as W_i , in an arbitrary ordering of attribute names; the Enron email data set [16] with more than 1.5 million documents (email messages and attachments) where all words, including attachments and envelope information, have been indexed resulting in about 1.2 million distinct words; and the ClueWeb09 [29] collection of crawled web-pages from which we extracted several databases of increasing size where the largest one was based on 0.4TB of HTML files with almost 3 billion of per-document-distinct words. The results of these tests show not only the suitability of our conjunction protocols for data sets of medium size (such as the Enron one) but demonstrate the scalability of these solutions to much larger databases (we target databases of one or two orders of magnitude larger). Existing solutions that are linear in the number of documents would be mostly impractical even for medium-size sets as the Enron case.

The pre-processing component accesses the three data sets as MySQL databases. The census data set is realized as a 100,000 record table, with one column for each attribute name. For the other two data sets, we wrote a set of Perl scripts to clean, stem, and convert the Enron emails and their attachments as well as the ClueWeb09 HTML files into separate MySQL tables. This transformation allow us to leverage the same pre-processing component for all data sets and it sets the foundation for handling more elaborate data sets such as XML document collections.

For the largest database derived from the ClueWeb09 data set based on 0.4 TB of HTML files, the sizes of the TSet hash table and XSet Bloom filter are 144.4 GB and 9.7 GB, respectively. The corresponding sizes for Enron and the census data are 12.4 GB and 0.6 GB, and 252 MB and 5 MB, respectively. For the multi-client configuration, in which each tuple includes the record decryption key encrypted with a key derived from the s-trap, the TSet size increases by 47%.

Experimental Results. The generation and retrieval of the DC (document cipher text) are straightforward and hence the measurements reported here cover only the protocols up to the discovery and decryption by the client of the rinds of the result set.

Figure 5 shows the latency of queries on one-term, called v , and three variants of two-term conjunctive queries on the Enron data set. In one-term queries, the selectivity of v varies from 3 to 690,492 documents. As this query consists only of an s-term, the figure illustrates that its execution time is linear in the cardinality of the corresponding TSet. The two-term conjunctive queries combine the previous queries with a fixed reference term. In the first of these queries, the fixed term acts as an x-term: each tuple retrieved from the TSet is checked against the XSet at the cost of an exponentiation. However, as we perform these operation in parallel to retrieving the TSet buckets from the disk, their cost is completely hidden by the disk I/O latency. Micro-benchmarks show that the average cost of retrieving a bucket which has a capacity of 10 tuples is comparable to $\sim 1,000$ single-threaded exponentiations. Similarly, the client-side exponentiation in the OXT protocol can be overlapped with disk and network I/O. It illustrates the for many surprising fact that even public-key cryptography, considered by many as expensive, can be surprisingly cheap when compared to the cost of accessing storage systems. The last two conjunctive queries use two fixed terms with different selectivity, α and β , as s-terms. Their invariable execution time is dominated by the cost of retrieving the TSet tuples corresponding to their s-terms, irrespective the variable selectivity of the xterm v : the two horizontal lines intersect with the single-term curve exactly where v corresponds to α and β , respectively. This illustrates the importance of s-term selection, as discussed in Section 3.1.1. Experiments on the census data show similar characteristics. All experiments were run on IBM Blades HS22 attached to a commodity SAN system.

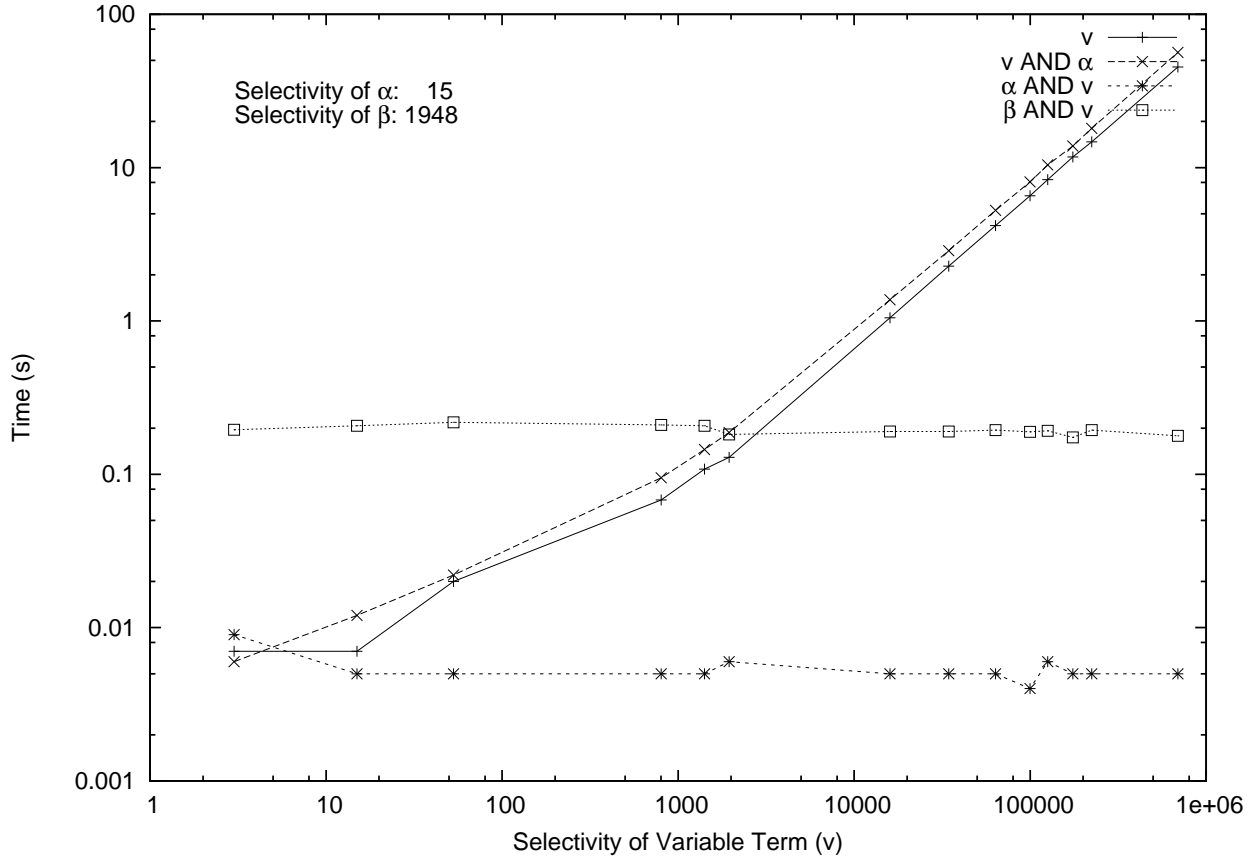


Figure 5: Enron Performance Measurement: Single Term & Conjunction

To further assess the scalability of our query engine, we generated several EDBs from increasingly larger subsets of the ClueWeb09 data set ranging from 408,450 to 13,284,801 HTML files having a size from 20 to 410 GBs and from 142,112,027 to 2,689,262,336 per-document-distinct words. To make these databases comparable, we injected some artificial and non-conflicting query terms to randomly selected documents simulating words of various selectivity. We then queried these terms in the various databases. Figure 6 shows the results and confirms that our implementation matches the scalability of our (theoretical) algorithms even when our databases exceed the size of available RAM: If the size of the result set is constant, then query time is largely independent of the size of the database and for result sets where the size is proportional to the database size, the cost is linear in the database size.

To quantify the ‘performance cost’ of the query and data privacy that our system provides, we measured the execution times of the same queries when run by a MySQL server (release 5.5) using its full text search capability on the Enron data set. To get a fair comparison, we measure execution times using a “cold” MySQL server, where we restarted the server and flushed the file-system caches before each query. In this scenario, our implementation is very competitive: on single term queries it is faster, sometimes significantly, for queries with small result sets and only slower (by a factor of 2) for terms with very poor selectivity, such as the very common keyword ‘enron’. For two-term conjunctions, the OXT prototype is slower than MySQL when the s-term is badly chosen. However, a good s-term selection results in consistently better performance than (cold) MySQL, in some cases by an order of magnitude.

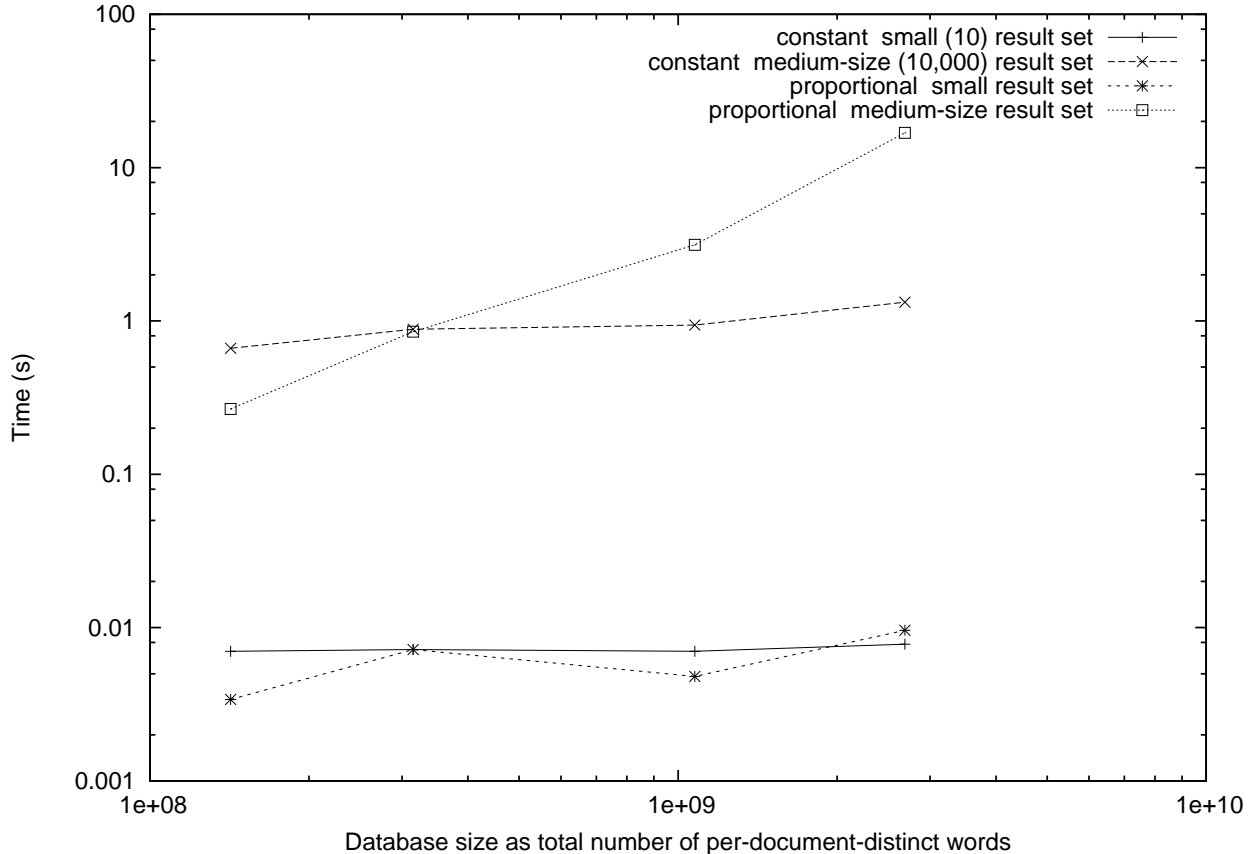


Figure 6: Clueweb09 Performance Measurement: Scaling Database Size

As the Enron data set is much smaller than our server’s RAM, a significant fraction of the data set and associated indexes get pre-fetched and cached in the MySQL and file system buffers after the first query. Not surprisingly, MySQL with warm caches is faster by one order of magnitude or more on most queries. We expect an alternative OXT prototype optimized for small data sets, namely using a RAM-resident TSet implementation, will be competitive with MySQL in warm cache scenarios. Finally, we need to emphasize that our prototype supports only a small subset of the query capabilities of a commercial database server, such as MySQL.

As our future work aims for even larger databases, we plan to investigate alternative strategies for realizing efficient disk-resident XSets.

8 Related Work

We discuss here the most directly relevant papers to our work [15, 14, 21, 3, 11]. For more on related work see [15, 14]. Most work on SSE has focused on single-keyword search. After several solutions with complexity linear in the number of documents, Curtmola et al. [15] present the first sublinear solution for single-keyword search whose complexity is linear in the number of matching documents (hence optimal). They also improve on previous security models, in particular by providing an adaptive security definition and solutions in this model. The SKS protocol from Figure 1 essentially follows their approach. Chase and Kamara [14] extend and generalize the security model to more complex data (e.g., graphs) and introduce the notion of associated data that allows to compose

different components of the protocol. We follow their model and modular approach. In addition, [14] introduce the notion of “controlled disclosure” that models the multi-client scenario mentioned above (see also the related “virtual private storage” setting of [25, 26]). Recently [28] introduced a strong security model of universal composability for searchable encryption, but their scheme supports only single-keyword search and it runs in time linear in the database size.

Conjunctions in the SSE setting were first considered by Golle, Staddon and Waters [21]. Their solutions are linear in the number d of all documents in the database, requiring $O(d)$ communication and same number of regular exponentiations at client and server, or communication that is only linear in the number of conjunctive terms but requires $O(d)$ pairing operations at the server. These solutions apply to structured (attribute-value) data only and leak the attributes being searched. On the positive side, they avoid leakage of keyword repetitions or access patterns other than the ciphertexts of matching documents. This work was followed up by Ballard et al. [3] and Byun et al. [11] but the above restrictions (linear in d complexity and application to structured-only data) remain. It needs to be noted that while our solutions have some crucial advantages in terms of performance and generality (applying to arbitrary data), they pay a price on privacy allowing significantly more leakage than [21, 3, 11]; on the other hand, a privacy advantage of our solutions is that they don’t disclose the searched attributes. An alternative SSE solution to conjunctions, also linear in the number of documents, can be obtained from [34, 30]; it uses per-document Bloom filters to decide for each document whether it contains all keywords in the conjunction.

CryptDB [33] supports most of SQL, including general boolean queries, over encrypted data and hence considerably more functionality than our work. It offers, though, also less privacy as CryptDB leaks statistical information on the *whole* queried column, not only related to matching rows, and this leakage can reveal significant information such as repeated values in the column (due to the use of deterministic encryption). Performance-wise, the schemes are hard to compare: All provided benchmarks in [33] were on databases which are several orders of magnitude smaller than ours and easily fit into RAM. Having small databases fitting into RAM side-steps the considerable challenge of efficiently dealing with high-latency storage, unavoidable for truly large databases.¹⁰ A related issue arises also with protocols based on the non-adaptive version of [15] which relies on a sequential linked-list implementation for the inverted index: the required inherently sequential list traversal will cause considerably latency on queries and limit scalability as soon as indices exceed the available RAM. Our implementation allows parallel list access which can hide these latency costs in overlapping network and computation costs to a mostly constant rather than linear factor. Lastly, we like to point out that our scheme could replace CryptDB’s use of the linear SSE scheme [36] for full-text search and, with corresponding support in the query optimizer, for less leaky boolean (sub-)expressions on single tables.

Another line of work on searchable encryption addresses support for efficient updates of the encrypted database [38, 27], but so far only single-keyword search was considered in this setting. Our OXT protocol, and its boolean query capabilities, can be extended to support efficient database updates at the cost of increased EDB size and partial leakage that can help differentiating between old and new documents [12].

Searchable encryption has also been investigated in the public-key domain [9, 39, 1, 4, 10, 35] allowing parties other than the data owner to encrypt into the database, but the cost of public key solutions makes the practical application of these techniques more limited (e.g., to a relatively small selection of tags or keywords). The case of conjunctive queries in the PK setting has been studied in [10].

¹⁰Recent work [37] extends CryptDB and demonstrates good performance also with bigger datasets, although still considerably smaller than the ones used in our own benchmarks and follow-up work [12].

9 Conclusions and Research Questions

The premise of this work is that in order to provide truly practical SSE solutions one needs to accept a certain level of leakage; therefore, the aim is to achieve an acceptable balance between leakage and performance, with formal analysis ensuring upper bounds on such leakage. Our solutions strike such a practical balance by offering performance that scales to very large data bases; supporting search in both structured and textual data with general Boolean queries; and confining leakage to access (to encrypted data) patterns and some query-term repetition only, with formal analysis defining and proving the exact boundaries of leakage. We stress that while in our solutions leakage never occurs in the form of direct exposure of plain data or searched values, when combined with side-information that the server may have (e.g., what are the most common searched words), such leakage can allow for statistical inference on plaintext data. Nonetheless, it appears that in many practical settings the benefits of search would outweigh moderate leakage (especially given the alternatives of outsourcing the plaintext data or keeping it encrypted but without the ability to run useful searches).

Our report on the characteristics and performance of our prototype points to the fact that scalability can only be achieved by low-complexity protocols which admit highly parallelizable implementations of their computational and I/O paths. Our protocols are designed to fulfill these crucial performance requirements.

There are interesting design and research challenges arising from this work. What we call “the challenge of being imperfect” calls for trade-offs between privacy and performance that can only be evaluated on the basis of a formal treatment of leakage. Understanding the limits of what is possible in this domain and providing formal lower bounds on such trade-offs appears as a non-trivial problem that deserves more attention. Some of these problems may still be unresolved even for plaintext data. The seemingly inherent difference pointed out in the introduction between the complexity of resolving conjunctions with high-frequency terms versus conjunctions with low-frequency terms, but with a similar-size result set, may be such a case. We do not know of a proven lower bound in this case although the work of Patrascu [32], for example, may point to some relevant conjectured bounds.

Another important evaluation of leakage is what we refer to as “semantic leakage.” How much can an attacker learn from the data given the formal leakage profile and side knowledge on the plaintext data? Clearly, the answer to this question is application-dependent but one may hope for some general theory in which these questions can be studied. The success of differential privacy in related domains opens some room for optimism in this direction. Demonstrating specific attacks in real-world settings is also an important direction to pursue. We note that in some settings just revealing the size of the number of documents matching a query may leak important information on the query contents (e.g., [23]). Therefore, developing masking techniques that include dummy or controlled data to obscure statistical information available to the attacker seems as an important research direction to strengthen the privacy of solutions as those developed here. ORAM-related techniques can be certainly help in this setting, especially given the progress on the practicality of these techniques in last years.

Yet another research direction is to characterize plaintext-search algorithms that lend themselves for adaptation to the encrypted setting. The s-term and x-term based search that we use is such an example: It treats data in “black-box” form that translates well to the encrypted setting. In contrast, search that looks at the data itself (e.g., sorting it) may not work in this setting or incur in significantly increased leakage (e.g., requiring order-preserving or deterministic encryption). Finally, it would be interesting to see more examples (in other two-party, or multi-party, protocols) of our approach, which is central to the design of OXT, of removing interaction from protocols by

pre-computing and storing some of the protocol messages during a pre-computation phase.

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. Abdalla, M. Bellare, D. Catalano, E. Kiltz, T. Kohno, T. Lange, J. Malone-Lee, G. Neven, P. Paillier, and H. Shi. Searchable encryption revisited: Consistency properties, relation to anonymous IBE, and extensions. In V. Shoup, editor, *CRYPTO 2005*, volume 3621 of *LNCS*, pages 205–222. Springer, Aug. 2005.
- [2] N. Alon and J. H. Spencer. *The Probabilistic Method*. Wiley, New York, 1992.
- [3] L. Ballard, S. Kamara, and F. Monrose. Achieving efficient conjunctive keyword searches over encrypted data. In S. Qing, W. Mao, J. López, and G. Wang, editors, *ICICS 05*, volume 3783 of *LNCS*, pages 414–426. Springer, Dec. 2005.
- [4] M. Bellare, A. Boldyreva, and A. O’Neill. Deterministic and efficiently searchable encryption. In A. Menezes, editor, *CRYPTO 2007*, volume 4622 of *LNCS*, pages 535–552. Springer, Aug. 2007.
- [5] M. Bellare, P. Rogaway, and T. Spies. The FFX mode of operation for format-preserving encryption. Draft 1.1, NIST Block Cipher Modes Working-Group, 2010.
- [6] B. H. Bloom. Space/time trade-offs in hash coding with allowable errors. *Communications of the Association for Computing Machinery*, 13(7):422–426, 1970.
- [7] D. Boneh and X. Boyen. Efficient selective-ID secure identity based encryption without random oracles. In C. Cachin and J. Camenisch, editors, *EUROCRYPT 2004*, volume 3027 of *LNCS*, pages 223–238. Springer, May 2004.
- [8] D. Boneh and X. Boyen. Short signatures without random oracles. In C. Cachin and J. Camenisch, editors, *EUROCRYPT 2004*, volume 3027 of *LNCS*, pages 56–73. Springer, May 2004.
- [9] 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. Springer, May 2004.
- [10] D. Boneh and B. Waters. Conjunctive, subset, and range queries on encrypted data. In S. P. Vadhan, editor, *TCC 2007*, volume 4392 of *LNCS*, pages 535–554. Springer, Feb. 2007.
- [11] J. W. Byun, D. H. Lee, and J. Lim. Efficient conjunctive keyword search on encrypted data storage system. In *EuroPKI*, pages 184–196, 2006.
- [12] D. Cash, J. Jagger, S. Jarecki, C. Jutla, H. Krawczyk, M. C. Rosu, and M. Steiner. Dynamic searchable encryption in very-large databases: Data structures and implementation. Manuscript, 2013.
- [13] 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. Springer, June 2005.
- [14] M. Chase and S. Kamara. Structured encryption and controlled disclosure. In *ASIACRYPT 2010*, *LNCS*, pages 577–594. Springer, Dec. 2010.

- [15] 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. ACM Press, Oct. 2006.
- [16] EDRM (edrm.net). Enron data set. <http://www.edrm.net/resources/data-sets/edrm-enron-email-data-set-v2>.
- [17] S. D. Galbraith, K. G. Paterson, and N. P. Smart. Pairings for cryptographers. *Discrete Applied Mathematics*, 156(16):3113–3121, 2008.
- [18] C. Gentry. Fully homomorphic encryption using ideal lattices. In M. Mitzenmacher, editor, *41st ACM STOC*, pages 169–178. ACM Press, May / June 2009.
- [19] E.-J. Goh. Secure indexes. Cryptology ePrint Archive, Report 2003/216, 2003. <http://eprint.iacr.org/>.
- [20] S. Goldwasser and R. Ostrovsky. Invariant signatures and non-interactive zero-knowledge proofs are equivalent (extended abstract). In E. F. Brickell, editor, *CRYPTO'92*, volume 740 of *LNCS*, pages 228–245. Springer, Aug. 1993.
- [21] P. Golle, J. Staddon, and B. R. Waters. Secure conjunctive keyword search over encrypted data. In M. Jakobsson, M. Yung, and J. Zhou, editors, *ACNS 04*, volume 3089 of *LNCS*, pages 31–45. Springer, June 2004.
- [22] IARPA. Security and Privacy Assurance Research (SPAR) Program - BAA, 2011. http://www.iarpa.gov/solicitations_spar.html/.
- [23] 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.
- [24] S. Jarecki, C. Jutla, H. Krawczyk, M. C. Rosu, and M. Steiner. Outsourced symmetric private information retrieval. 20th ACM Conference on Computer and Communications Security, CCS 2013, November 2013, Berlin, Germany.
- [25] S. Kamara and K. Lauter. Cryptographic cloud storage. In *Financial Cryptography Workshops*, pages 136–149, 2010.
- [26] S. Kamara, C. Papamanthou, and T. Röder. "CS2: A searchable cryptographic cloud storage system", 2011. <http://research.microsoft.com/pubs/148632/CS2.pdf>.
- [27] S. Kamara, C. Papamanthou, and T. Roeder. Dynamic searchable symmetric encryption. In *Proc. of CCS'2012*, 2012.
- [28] K. Kurosawa and Y. Ohtaki. Uc-secure searchable symmetric encryption. In *Financial Cryptography*, page 285, 2012.
- [29] Lemur Project. ClueWeb09 dataset. <http://lemurproject.org/clueweb09.php/>.
- [30] V. Pappas, M. Raykova, B. Vo, S. M. Bellovin, and T. Malkin. Private search in the real world. In *Annual Computer Security Applications Conference*, 2011.
- [31] V. Pappas, B. Vo, F. Krell, S. G. Choi, V. Kolesnikov, A. Keromytis, and T. Malkin. Blind Seer: A Scalable Private DBMS. Manuscript, 2013.
- [32] M. Patrascu. Towards polynomial lower bounds for dynamic problems. In *42nd ACM STOC*, pages 603–610. ACM Press, 2010.
- [33] R. A. Popa, C. M. S. Redfield, N. Zeldovich, and H. Balakrishnan. CryptDB: Protecting confidentiality with encrypted query processing. In *Proceedings of the 23rd ACM Symposium on Operating Systems Principles (SOSP'11)*. ACM, Oct. 2011.
- [34] M. Raykova, B. Vo, T. Malkin, and S. M. Bellovin. Secure anonymous database search. In *ACM Cloud Computing Security Workshop*, 2009.

- [35] E. Shi, J. Bethencourt, H. T.-H. Chan, D. X. Song, and A. Perrig. Multi-dimensional range query over encrypted data. In *2007 IEEE Symposium on Security and Privacy*, pages 350–364. IEEE Computer Society Press, May 2007.
- [36] D. X. Song, D. Wagner, and A. Perrig. Practical techniques for searches on encrypted data. In *2000 IEEE Symposium on Security and Privacy*, pages 44–55. IEEE Computer Society Press, May 2000.
- [37] S. Tu, M. F. Kaashoek, S. Madden, and N. Zeldovich. Processing analytical queries over encrypted data. In *Proceedings of the 39th International Conference on Very Large Data Bases (VLDB)*, 2013.
- [38] 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.
- [39] B. R. Waters, D. Balfanz, G. Durfee, and D. K. Smetters. Building an encrypted and searchable audit log. In *NDSS 2004*. The Internet Society, Feb. 2004.

A Proof of Lemma 4

Let A be an adversary and Π be the OXT protocol, and let $G_0 = \mathbf{Cor}_A^\Pi(\lambda)$. We want to show

$$\Pr[G_0 = 1] = \text{neg}(\lambda).$$

We first modify G_0 to output 1 if the T-set ever errs in returning the correct data \mathbf{t} when the server is processing a search in the game. It is straightforward to construct an adversary B' such that

$$\Pr[G_1 = 1] - \Pr[G_0 = 1] = \mathbf{AdvCor}_{B'}^\Pi(\lambda).$$

Next we consider a modification of G_0 that starts by choosing random functions f_X, f_I mapping $\{0, 1\}^*$ to Z_p^* , and then replaces all evaluations of $F_p(K_X, \cdot)$ and $F_p(K_I, \cdot)$ with evaluations of $f_X(\cdot)$ and $f_I(\cdot)$ respectively.

Call this modification G_2 . By a straightforward reduction the PRF security of F_p and a hybrid argument, we can build an efficient adversary B such that

$$\Pr[G_2 = 1] - \Pr[G_1 = 1] \leq 2 \cdot \mathbf{Adv}_{F_p, B}^{\text{prf}}(\lambda).$$

We complete the proof by showing that

$$\Pr[G_2 = 1] \leq N^2/(p-1) + (d+m)/p$$

and combining and rearranging the inequalities.

In G_2 , let $\bar{w} = (w_1, w_2)$ be the query chosen by A after seeing EDB. The game will output 1 only if the simulated search protocol causes the client to have output not equal to the permuted results in $\text{DB}(\bar{w})$. Let us examine the search protocol computation in more detail. By the correctness of the T-set implementation, the vector \mathbf{t} computed by the server consists of ciphertexts decrypting exactly to the set $\text{DB}(w_1)$. Of these, the server returns those that satisfy $g^{f_I(\text{ind}) \cdot f_X(w_2)} \in \text{XSet}$. If $\text{ind} \in \text{DB}(w_2)$ then this will clearly be the case, so we only need to rule out the possibility that a false positive occurs and the test with the XSet is true but $\text{ind} \notin \text{DB}(w_2)$ for this to happen, there would have to be some $(\text{ind}', w') \neq (\text{ind}, w_2)$ such that

$$g^{f_I(\text{ind}) \cdot f_X(w_2)} = g^{f_I(\text{ind}') \cdot f_X(w')}.$$

Assume that for every ind , $f_I(\text{ind}) \neq 0$ and for every keyword w , $f_X(w) \neq 0$. (This is true with all but $(d+m)/p$ probability.) Then the above question holds with probability $1/(p-1)$ over the choice of f_I and f_X . The desired bound follows from taking a union bound over all N^2 possibilities for (ind, w_2) and (ind', w') , and by adding the probability \mathbf{AdvCor}_A^Π that in the SSE correctness game with adversary A there occurs a correctness break in the underlying T-set implementation Π .

B Proof of Theorem 5

Here we prove Theorem 5. More precisely, we prove a more general version of it that considers a generic instantiation of TSet. See Theorem 9 below.

We recall the non-adaptive security definition used in the theorem.

Definition 8 Let $\Pi = (\text{EDBSetup}, \text{Search})$ be an SSE scheme and let \mathcal{L} be an algorithm. For efficient algorithms A and S , we define experiments (algorithms) $\mathbf{Real}_A^\Pi(\lambda)$ and $\mathbf{Ideal}_{A,S}^\Pi(\lambda)$ as follows:

Real $_A^\Pi(\lambda)$: $A(1^\lambda)$ chooses DB and a list of queries \mathbf{q} . The experiment then runs $(K, \text{EDB}) \leftarrow \text{EDBSetup}(\text{DB})$. For each $i \in |\mathbf{q}|$, it runs the Search protocol with client input $(K, \mathbf{q}[i])$ and server input EDB and stores the transcript and the client's output in $\mathbf{t}[i]$. Finally the game gives EDB and \mathbf{t} to A , which returns a bit that the game uses as its own output.

Ideal $_{A,S}^\Pi(\lambda)$: $A(1^\lambda)$ chooses DB and a list of queries \mathbf{q} . The experiment then runs $S(\mathcal{L}(\text{DB}, \mathbf{q}))$ and gives its output to A , which returns a bit that the game uses as its own output.

We say that Π is \mathcal{L} -semantically-secure against non-adaptive attacks if for all efficient 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 describe the leakage function \mathcal{L} that uses \mathcal{L}_{ext} along with the T-set leakage function \mathcal{L}_T . On input DB and (\mathbf{s}, \mathbf{x}) , it computes the vector \mathbf{T} via

```

For  $w \in W$  do
   $K \xleftarrow{\$} \{0, 1\}^\lambda$ ;  $\mathbf{t} \leftarrow \perp$ 
  For  $c = 1, \dots, T_w$  do  $y \xleftarrow{\$} Z_p^*$ ;  $\mathbf{e} \leftarrow \text{Enc}(K, 0^\lambda)$ ;  $\mathbf{t}[c] \leftarrow (y, \mathbf{e})$ 
   $\mathbf{T}[w] \leftarrow \mathbf{t}$ 
End.
```

Then it outputs $(\mathcal{L}_{\text{ext}}(\text{DB}, (\mathbf{s}, \mathbf{x})), \mathcal{L}_T(\mathbf{T}, \mathbf{s}), \mathbf{T}[\mathbf{s}])$, where the notation $\mathbf{T}[\mathbf{s}] = (\mathbf{T}[\mathbf{s}[1]], \dots, \mathbf{T}[\mathbf{s}[Q]])$.

Theorem 9 Let \mathcal{L} be the leakage function defined above. Then SSE scheme OXT is \mathcal{L} -semantically-secure against non-adaptive attacks where all queries are 2-conjunctions, assuming that the DDH assumption holds in G , that F and F_p are secure PRFs, that (Enc, Dec) is an IND-CPA secure symmetric encryption scheme, and that Σ is a (non-adaptively) \mathcal{L}_T -secure and computationally correct T-set instantiation.

Proof. We structure our proof using several games G_0, G_1, \dots . In each game, A starts by supplying DB, \mathbf{q} , which is then given to an Initialize routine that produces an output, which is given to A who then outputs a bit that becomes the game output. Game G_0 is designed to generate exactly the same distribution as $\mathbf{Real}_A^\Pi(\lambda)$ (assuming no false positives occur) and the final game is structured so that it is easy to simulate exactly given the leakage profile instead of the actual DB, \mathbf{q} input. By relating the games, we can argue that the final simulator satisfies Definition 8 with OXT, completing the proof. For simplicity only, these games all model a version of the protocol where the client sends the maximum-sized xtrap (say, T group elements where T is some publicly known upper bound) on each query, and the proof is easily generalizable to the version where the server interactively tells the client to stop.

Game G_0 . The first game G_0 , which uses the routines in Figure 7, is an implementation of the real game with some minor changes that will make the analysis easier later. The game starts by running Initialize, which passes $(\text{DB}, \mathbf{s}, \mathbf{x})$ from A and simulates EDBSetup(DB), with a minor bookkeeping changes. In particular, while building the \mathbf{T} it records the permutations σ in a vector WPerms indexed by keywords.

<pre> Initialize(DB, s, x) // G₀ K_S, K_I, K_X, K_Z \xleftarrow{s} {0, 1}^λ (ind_i, W_i)_{i=1}^d \leftarrow DB For w ∈ W (ind₁, ..., ind_{T_w}) \leftarrow DB(w) σ \xleftarrow{s} Perm([T_w]); WPerms[w] \leftarrow σ K_e \leftarrow F(K_S, w); t \leftarrow ⊥ For c = 1, ..., T_w do e \leftarrow Enc(K_e, ind_{σ(c)}); xind \leftarrow F_p(K_I, ind_{σ(c)}) z \leftarrow F_p(K_Z, w c); y \leftarrow xind · z⁻¹ mod p t[c] \leftarrow (y, e) End T[w] \leftarrow t End (TSet, K_T) \leftarrow TSetSetup(T) For i = 1, ..., Q do STags[i] \leftarrow TSetGetTag(K_T, s[i]) XSet \leftarrow XSetSetup(K_X, K_I, DB) EDB \leftarrow (TSet, XSet) For i = 1, ..., Q do t[i] \leftarrow GenTrans(K_X, K_Z, s[i], x[i], STags[i]) Return (EDB, t) </pre>	<pre> XSetSetup(K_X, K_I, DB) // G₀ (ind_i, W_i)_{i=1}^d \leftarrow DB; XSet \leftarrow ∅ For each w ∈ W and ind ∈ DB(w) do e \leftarrow F_p(K_X, w) xind \leftarrow F_p(K_I, ind); h \leftarrow g^{e·xind} XSet \leftarrow XSet ∪ {h} End Return XSet GenTrans(EDB, K_X, K_Z, s, x, stag) // G₀ e \leftarrow F_p(K_X, x) For c = 1, ..., T do z \leftarrow F_p(K_Z, s c); xtoken[c] \leftarrow g^{e·z} Res \leftarrow ServerSearch(EDB, (stag, xtoken)) ResInds \leftarrow DB(s) ∩ DB(x) Return ((stag, xtoken), Res, ResInds) </pre>
--	---

Figure 7: Game G_0 .

Before generating the transcripts, the game computes an array **STags** of all of the **stag** values used in the game. For each $i = 1, \dots, Q$ it lets $\text{STags}[i] \leftarrow \text{TSetGetTag}(K_T, s[i])$. To compute the transcript array **t**, for $i = 1, \dots, Q$ it sets $t[i]$ to the output of $\text{GenTrans}(\text{EDB}, K_X, K_Z, s[i], x[i], \text{STags}[i])$, which is defined in the figure. There we use subroutine **ServerSearch** which we take to be the server's computation defined in **OXT** in response to the first client message.

The routine **GenTrans** generates a transcript as in the real game, except that it computes the **ResInds** array differently: Instead of decrypting the ciphertexts returned with **Res**, it looks up the **ind** values that correspond to the results (specifically, it computes $\text{DB}(s[i])$ and then finds the **ind** values amongst them that are also in $\text{DB}(x[i])$).

By design, G_0 is exactly $\text{Real}_A^\Pi(\lambda)$, except that false positives are assumed to never happen. By Theorem 4, assuming that F_p is a secure PRF, we have

$$\Pr[G_0 = 1] \leq \Pr[\text{Real}_A^\Pi(\lambda) = 1] + \text{neg}(\lambda).$$

Game G_1 . In the next game, G_1 , is exactly the same except we replace every evaluation of $F(K_S, \cdot), F_p(K_X, \cdot), F_p(K_I, \cdot), F_p(K_Z, \cdot)$ with evaluations of independent random functions with the appropriate domain and range. It is described in Figure 8. Note that since $F(K_S, \cdot)$ is never evaluated on the same input twice, we can equivalently replace its evaluations with random selections from the range (this is done in **Initialize** when selecting K_e). By a standard hybrid argument it is easy to show that there exist efficient adversaries $B_{1,1}, B_{1,2}, B_{1,3}$ such that

$$\Pr[G_1 = 1] - \Pr[G_0 = 1] \leq \text{Adv}_{F, B_{1,1}}^{\text{prf}}(\lambda) + 3 \cdot \text{Adv}_{F_p, B_{1,2}}^{\text{prf}}(\lambda) + \text{Adv}_{F_p, B_{1,3}}^{\text{prp}}(\lambda)$$

We omit the tedious and standard details of the adversaries. The salient feature of the game is that these keys are chosen at random, and then afterwards only used as key inputs to the corresponding functions.

<pre> Initialize(DB, s, x) // G₁, G₂ f_I, f_X, f_Z $\xleftarrow{\\$}$ Fun({0, 1}^λ, Z_p[*]) (ind_i, W_i)_{i=1}^d \leftarrow DB For w ∈ W (ind₁, ..., ind_{T_w}) \leftarrow DB(w) σ $\xleftarrow{\\$}$ Perm([T_w]); WPerms[w] \leftarrow σ K_e $\xleftarrow{\\$}$ {0, 1}^λ; t \leftarrow ⊥ For c = 1, ..., T_w do e \leftarrow Enc(K_e, ind_{σ(c)}) e \leftarrow Enc(K_e, 0^λ) xind \leftarrow f_I(ind_{σ(c)}) z \leftarrow f_Z(w c) y \leftarrow xind · z⁻¹ mod p t[c] \leftarrow (y, e) End T[w] \leftarrow t End (TSet, K_T) \leftarrow TSetSetup(T) For i = 1, ..., Q do STags[i] \leftarrow TSetGetTag(K_T, s[i]) XSet \leftarrow XSetSetup(f_X, f_I, DB) EDB \leftarrow (TSet, XSet) For i = 1, ..., Q do t[i] \leftarrow GenTrans(f_X, f_Z, s[i], x[i], STags[i]) Return (EDB, t) </pre>	<pre> XSetSetup(f_X, f_I, DB) // G₁, G₂ (ind_i, W_i)_{i=1}^d \leftarrow DB XSet \leftarrow ∅ For each w ∈ W and ind ∈ DB(w) do e \leftarrow f_X(w); xind \leftarrow f_I(ind); h \leftarrow g^{e·xind} XSet \leftarrow XSet ∪ {h} End Return XSet GenTrans(EDB, f_X, f_Z, s, x, stag) // G₁, G₂ e \leftarrow f_X(x) For c = 1, ..., T do z \leftarrow f_Z(s c); xtoken[c] \leftarrow g^{e·z} Res \leftarrow ServerSearch(EDB, (stag, xtoken)) ResInds \leftarrow DB(s) ∩ DB(x) Return ((stag, xtoken), Res, ResInds) </pre>
---	---

Figure 8: Games G_1 and G_2 . G_2 includes the boxed code and G_1 does not.

Game G_2 . In G_2 we modify G_1 to include the boxed code. This means that, during Initialize, the ciphertext e is always overwritten with an encryption of 0^λ (under the same key). We claim that there exists an efficient adversary B_2 such that

$$\Pr[G_2 = 1] - \Pr[G_1 = 1] \leq m \cdot \mathbf{Adv}_{\Sigma, B_2}^{\text{ind-cpa}}(\lambda)$$

This follows by a standard hybrid argument over the m encryption keys used in building \mathbf{T} for TSetSetup. We omit the tedious details. We stress that the reduction is possible because the game never invokes the decryption algorithm Dec of the scheme, meaning that the reduction does not need to decrypt ciphertexts during the IND-CPA game.

Game G_3 . In game G_3 (Figure 9) we alter the way the game is computed without changing its distribution. Intuitively, it precomputes the values in the XSet as well as all of the group elements that will be tested against the XSet, along with the group elements used xtoken arrays in transcripts that do not correspond to possible matches.

In Initialize, two arrays H, Y are filled in for use in XSetSetup and GenTranscript. H is indexed by an identifier (string in $\{0, 1\}^\lambda$) and a keyword from W , and holds group elements from G , and Y is indexed by a keyword and a number between 1 and T and also holds elements from G .

Now XSetSetup is modified to use the values from the H array. For a given $w \in W$ and $\text{ind} \in \text{DB}(w)$, it adds the value $H[\text{ind}, w]$ to XSet - This value is set to $g^{f_I(\text{ind}) \cdot f_X(w)}$ during Initialize, which is the same value that was computed in XSetSetup in game G_2 .

The H and Y arrays are both used in GenTranscript. We claim that, for any input DB and fixed

randomness, each invocation of `GenTranscript` will return exactly the same output in G_2 and G_3 . Clearly the `stag`, `Res`, `ResInds` values will be the same, so we focus on the `xtoken` array.

In G_2 , for each c , `GenTranscript`(`EDB`, f_X , f_Z , s , x , `stag`) sets `xtoken`[c] to $g^{f_X(x) \cdot f_Z(s \parallel c)}$. In G_3 , `GenTranscript`(`EDB`, f_X , f_Z , s , x , `stag`) looks up $\sigma = \text{WPerms}[s]$, $\text{DB}(s) = (\overline{\text{ind}}_1, \dots, \overline{\text{ind}}_{T_s})$ and \mathbf{t} . By the correctness of the `TSet` instantiation we have $\mathbf{t} = (f_I(\overline{\text{ind}}_{\sigma(c)})/f_Z(s \parallel c), \mathbf{e}_c)_{c=1}^{T_s}$, for some ciphertexts $\mathbf{e}_1, \dots, \mathbf{e}_{T_s}$, where the $\overline{\text{ind}}_c$ and σ are the same values as above.

Then for $c = 1, \dots, T_s$, y is set to $f_I(\overline{\text{ind}}_{\sigma(c)})/f_Z(s \parallel c)$ and `xtoken`[c] is set

$$H[\overline{\text{ind}}_{\sigma(c)}, x]^{1/y} = g^{f_X(x) \cdot f_Z(s \parallel c)},$$

which shows that the `xtoken`[c] value is the same in both games for $c \leq T_s$. Note that here we have crucially used the fact that the set $\text{DB}(s)$ was parsed in the same canonical order in both places in the game.

For $c > T_s$, G_3 sets `xtoken`[c] to

$$Y[s, x, c] = g^{f_X(x) \cdot f_Z(s \parallel c)},$$

which is exactly the same as `xtoken`[c] in G_2 . This completes the claim that same value is returned in either game, so we have

$$\Pr[G_3 = 1] = \Pr[G_2 = 1].$$

Game G_4 . Game G_4 is exactly like G_3 except that the boxed code is included: Now the game simply sets every y value of the \mathbf{t} arrays to be an independently chosen element of Z_p^* . We claim that

$$\Pr[G_4 = 1] = \Pr[G_3 = 1].$$

This can be seen by observing that the random function f_Z is chosen during `Initialize` and then never evaluated in any of the other subroutines of G_3 , thanks to the modifications made in moving to G_3 . Moreover, for any $w \in \mathbf{W}$ and $c = 1, \dots, T_w$, the value of $z = f_Z(w \parallel c)$ is used exactly once during `Initialize`, and this value is uniform and independent of the rest of the randomness in the game. Thus the value $f_I(\overline{\text{ind}}_{\sigma(c)})/f_Z(w \parallel c)$ is also uniform and independent of the rest of the game, which justifies the change to simply selecting y at random. We note that the random function f_Z is also evaluated in the loop below that computes Y , but this loop evaluates $f_Z(w \parallel c)$ for $c > T_w$, meaning it uses different points of f_Z .

Game G_5 . In G_5 , instead of computing the values of the H and Y arrays as before, they are all selected at random from G . We claim there exist an efficient adversary B_5 such that

$$\Pr[G_5 = 1] - \Pr[G_4 = 1] \leq \text{Adv}_{G, B_5}^{\text{ddh}}(\lambda)$$

Let A be an adversary running in G_4 or G_5 . Using A , we build an adversary B_5 for the DDH problem. By Lemma 3, it suffices to build an adversary B'_5 that solves the extended version of DDH defined in the lemma.

B'_5 takes as input $(g^{\mathbf{a}}, g^{\mathbf{b}}, \mathbf{M})$ where $g^{\mathbf{a}} \in G^m$, $g^{\mathbf{b}} \in G^{d+T \cdot m}$ and $\mathbf{M} \in G^{m \times (d+T \cdot m)}$. Intuitively, it will use the $g^{\mathbf{a}}$ vector in place of the X array in the games, and the entries from $g^{\mathbf{b}}$ will play two roles: The discrete logarithms of the first d entries will take the place of the `xind` values in the loop computing the H array, and the discrete logarithms of the rest of the values will take the place of the $f_Z(w \parallel c)$ values in the loop populating Y . This means that some entries of $g^{\mathbf{b}}$ and \mathbf{M} may not be used, depending on the structure of the DB input.

B'_5 starts by running A , which outputs $(\text{DB}, (s, \mathbf{x}))$. B then simulates G_5 , except for the following changes: It does not choose f_X and f_I , and the loop that populates the X and H arrays is instead

<pre> Initialize(DB, s, x) // $G_3, \boxed{G_4}, \boxed{\boxed{G_5}}$ $f_I, f_X, f_Z \xleftarrow{\\$} \text{Fun}(\{0, 1\}^\lambda, Z_p^*)$ $(\text{ind}_i, W_i)_{i=1}^d \leftarrow \text{DB}$ For $w \in W$ and each ind_i do $e \leftarrow f_X(w); X[w] \leftarrow g^e; \text{xind} \leftarrow f_I(\text{ind}_i)$ $H[\text{ind}_i, w] \leftarrow X[w]^{\text{xind}}$ $\boxed{H[\text{ind}_i, w] \xleftarrow{\\$} G}$ End For $w \in W$ $(\overline{\text{ind}}_1, \dots, \overline{\text{ind}}_{T_w}) \leftarrow \text{DB}(w)$ $K_e \xleftarrow{\\$} \{0, 1\}^\lambda; \mathbf{t} \leftarrow \perp$ $\sigma \xleftarrow{\\$} \text{Perm}([T_w]); \text{WPerms}[w] \leftarrow \sigma$ For $c = 1, \dots, T_w$ do $\text{xind} \leftarrow f_I(\overline{\text{ind}}_{\sigma(c)}); z \leftarrow f_Z(w \parallel c)$ $y \leftarrow \text{xind} \cdot z^{-1} \bmod p$ $\boxed{y \xleftarrow{\\$} Z_p^*}; e \leftarrow \text{Enc}(K_e, 0^\lambda)$ $\mathbf{t}[c] \leftarrow (y, e)$ End $\mathbf{T}[w] \leftarrow \mathbf{t}$ For $u \in W \setminus \{w\}$ do For $c = T_w + 1, \dots, T$ do $Y[w, u, c] \leftarrow X[u]^{f_Z(w \parallel c)}$ $\boxed{Y[w, u, c] \xleftarrow{\\$} G}$ End End End End End $(\text{TSet}, K_T) \leftarrow \text{TSetSetup}(\mathbf{T})$ For $i = 1, \dots, Q$ do $\text{STags}[i] \leftarrow \text{TSetGetTag}(K_T, \mathbf{s}[i])$ $\text{XSet} \leftarrow \text{XSetSetup}(\text{DB}, H)$ $\text{EDB} \leftarrow (\text{TSet}, \text{XSet})$ For $i = 1, \dots, Q$ do $\mathbf{t}[i] \leftarrow \text{GenTrans}(\text{DB}, \text{EDB}, H, \mathbf{s}[i], \mathbf{x}[i], \text{STags}[i])$ Return (EDB, \mathbf{t}) </pre>	<pre> XSetSetup(DB, H) // G_3, G_4, G_5 $(\text{ind}_i, W_i)_{i=1}^d \leftarrow \text{DB}$ $\text{XSet} \leftarrow \emptyset$ For each $w \in W$ and $\text{ind} \in \text{DB}(w)$ do $\text{XSet} \leftarrow \text{XSet} \cup \{H[\text{ind}, w]\}$ Return XSet GenTrans(DB, EDB, H, s, x, stag) // G_3, G_4, G_5 $\mathbf{t} \leftarrow \text{TSetRetrieve}(\text{TSet}, \text{stag})$ $(\overline{\text{ind}}_1, \dots, \overline{\text{ind}}_{T_s}) \leftarrow \text{DB}(s); \sigma \leftarrow \text{WPerms}[s]$ For $c = 1, \dots, T_s$ do $(y, e) \leftarrow \mathbf{t}[c]; \text{xtoken}[c] \leftarrow H[\text{ind}_{\sigma(c)}, x]^{1/y}$ For $c = T_s + 1, \dots, T$ do $\text{xtoken}[c] \leftarrow Y[s, x, c]$ Res $\leftarrow \text{ServerSearch}(\text{EDB}, (\text{stag}, \text{xtoken}))$ ResInds $\leftarrow \text{DB}(s) \cap \text{DB}(x)$ Return $((\text{stag}, \text{xtoken}), \text{Res}, \text{ResInds})$ </pre>
--	--

Figure 9: Games G_3, G_4 and G_5 . G_5 includes the doubly boxed code and singly boxed code, G_4 includes only the singly boxed code, and G_3 includes neither.

computed by taking the i -th keyword w (in arbitrary order), it places $g^{\mathbf{a}[i]}$ in $X[w]$, and for the ind_j , it places $\mathbf{M}[i, j]$ in $H[\text{ind}_i, w]$. Finally, when computing the Y array, for the entry at (w, u, c) , where w is the i -th keyword in the outer loop (in the same order as the H array loop) and u is the j -th keyword in the inner loop, it uses the value from $\mathbf{M}[i, d + j \cdot T + c]$.

It simulates the rest of the game exactly as specified, and finally outputs whatever A outputs. If the matrix \mathbf{M} is computed as $\mathbf{M} = g^{\mathbf{ab}^\top}$, then it is apparent that B'_5 simulates G_4 for A , and thus

$$\Pr[B'_5(g, g^{\mathbf{a}}, g^{\mathbf{b}}, g^{\mathbf{ab}^\top}) = 1] = \Pr[G_4 = 1].$$

<pre> Initialize(DB, s, x) // G₆, G₇, G₈ (ind_i, W_i)_{i=1}^d ← DB For w ∈ W and i ∈ [d] do H[ind_i, w] $\xleftarrow{\\$}$ G For w ∈ s do WPerms[w] $\xleftarrow{\\$}$ Perm([T_s]) For w ∈ W do K $\xleftarrow{\\$}$ {0, 1}^λ; t ← ⊥ For c = 1, ..., T_w do y $\xleftarrow{\\$}$ Z_p[*]; e ← Enc(K, 0^λ); t[c] ← (y, e) End T[w] ← t End (TSet, STags) ← S_T(L_T(DB, s), T[s]) XSet ← XSetSetup(DB, H) EDB ← (TSet, XSet) For i = 1, ..., Q do t[i] ← GenTrans(DB, EDB, H, s[i], x[i], STags[i]) End Return (EDB, t) XSetSetup(DB, H) // G₇, G₈ XSet ← ∅ For w ∈ W and ind ∈ DB(w) do If ∃ i : ind ∈ DB(s[i]) ∧ x[i] = w then XSet ← XSet ∪ {H[ind, w]} Else h $\xleftarrow{\\$}$ G; XSet ← XSet ∪ {h} End Return XSet </pre>	<pre> GenTrans(DB, EDB, H, s, x, stag, i) // G₈ t ← TSetRetrieve(TSet, stag) (ind₁, ..., ind_{T_s}) ← DB(s); σ ← WPerms[s] For c = 1, ..., T_s do (y, e) ← t[c] If ind_{σ(c)} ∈ DB(s) ∩ DB(x) then xtoken[c] ← H[ind_{σ(c)}, x]^{1/y} Else If ∃ j ≠ i : ind_{σ(c)} ∈ DB(s[j]) ∧ x[j] = x then xtoken[c] ← H[ind_{σ(c)}, x]^{1/y} Else xtoken[c] $\xleftarrow{\\$}$ G End For c = T_s + 1, ..., T do xtoken[c] $\xleftarrow{\\$}$ G Res ← ServerSearch(EDB, (stag, xtoken)) ResInds ← DB(s) ∩ DB(x) Return ((stag, xtoken), Res, ResInds) </pre>
---	--

Figure 10: Games G_6, G_7, G_8 .

If instead \mathbf{M} is selected at random from $G^{m \times (d+mT)}$, then B'_5 simulates G_5 , so

$$\Pr[B'_5(g, g^{\mathbf{a}}, g^{\mathbf{b}}, \mathbf{M}) = 1] = \Pr[G_5 = 1].$$

Subtracting these inequalities shows that B'_5 satisfies the desired relation.

Game G_6 . The Initialize code for G_6 is described in Figure 10, and its other routines remain as defined in G_5 . G_6 deletes some irrelevant code (like selecting the random functions) and generates the TSet using a simulator which we will show to be guaranteed to exist by the TSet security notion.

Before describing G_6 in detail, we consider the following TSet adversary B_6 . In the T-set games, B_6 starts by generating $(\text{DB}, \mathbf{s}, \mathbf{T})$ exactly as described in G_5 . Upon receiving $(\text{TSet}, \text{STags})$ from its game, B_6 continues to simulate the game with these values until A halts, and it outputs whatever A output. By the \mathcal{L}_T -security of the T-set instantiation, there exists an efficient simulator \mathcal{S}_T for B_6 .

In game G_6 , we replace the call to TSetSetup and loop generating STags with a call to the simulator \mathcal{S}_T on input $(\mathcal{L}_T(\text{DB}, \mathbf{s}), \mathbf{T}[\mathbf{s}])$. By the construction of G_5 and G_6 , we have that the T-set real game with adversary B_6 perfectly simulates G_5 , while the T-set ideal game with B_6 perfectly simulates G_6 . Thus a non-negligible difference in $\Pr[G_6 = 1]$ and $\Pr[G_5 = 1]$ would contradict the assumed \mathcal{L}_T -security of the T-set instantiation, giving

$$\Pr[G_6 = 1] - \Pr[G_5 = 1] \leq \text{neg}(\lambda).$$

Games G_7 and G_8 . Games G_7 and G_8 change the way the H array is accessed in order to enable the final simulator to work with its given leakage. Intuitively, now whenever the game access the H array at an index (ind, x) , it first tests to see if the game will ever access that index in H again. If it will come back to this position, it uses the value from H . If not, then the game replaces the H access with a random choice. Since that was to be the only usage of that position of H during the game, this doesn't affect the distribution of the game.

We now describe how this is implemented. G_7 only changes the way **XSet** works. Recall that after H is populated, it is used in two places: In **XSetSetup** and in **GenTranscript**. The routine **XSetSetup** will never repeat an access to H , so for an index (ind, w) it only needs to test if **GenTranscript** will read this position. But **GenTranscript** will only read positions such that $\text{ind} \in \text{DB}(\mathbf{s}[i])$ and $x = \mathbf{x}[i]$ for some i , and this is exactly what **XSetSetup** tests for in G_7 . By this observation and the discussion above, we have

$$\Pr[G_7 = 1] = \Pr[G_6 = 1].$$

For G_8 we change the way **GenTranscript** accesses H . To detect a possible repeated access, it must test if either **XSetSetup** will read that index, or if **GenTranscript** will read it again. In G_7 we modified **XSetSetup** so that it only accesses a position $H[\text{ind}, w]$ if $\text{ind} \in \text{DB}(\mathbf{s}[i]) \cap \text{DB}(\mathbf{x}[i])$, so the first “If” statement in the G_8 version of **GenTranscript** tests for this. What remains to calculate are the positions which **GenTranscript** will access twice. First, we have that a repeated position must occur during two different calls to **GenTranscript** because one execution of the subroutine only touches unique indices of H . So for an index (ind, x) to be accessed twice, we must have that ind is a member of both $\text{DB}(\mathbf{s}[i])$ and $\text{DB}(\mathbf{s}[j])$, for some $i \neq j$. For the x to repeat, we must have $\mathbf{x}[i] = \mathbf{x}[j]$. This condition for a repeated access is exactly what the “Else If” statement in **GenTranscript** tests for. An argument similar to the one for the previous game transition gives

$$\Pr[G_8 = 1] = \Pr[G_7 = 1].$$

Simulator. We next give a simulator \mathcal{S} that uses its leakage to compute the same distribution as G_8 . By collecting the relations between the games, we will show that \mathcal{S} satisfies the requirement in the theorem.

\mathcal{S} takes as input $\mathcal{L}(\text{DB}, \mathbf{s}, \mathbf{x})$, which consists of $(N, \bar{\mathbf{s}}, \text{SP}, \text{RP}, \text{IP}, \mathcal{L}_T(\text{DB}, \mathbf{s}), \mathbf{T}[\mathbf{s}])$ and must output a simulated $\text{EDB} = (\text{TSet}, \text{XSet})$ and transcript array \mathbf{t} .

Our simulator will first compute what we call the *restricted equality pattern* of \mathbf{x} , denoted $\hat{\mathbf{x}}$. Intuitively, $\hat{\mathbf{x}}$ describes which \mathbf{x} -terms are “known” to be equal by the server.

The restricted equality pattern $\hat{\mathbf{x}} \in [m]^Q$ is computed by \mathcal{S} as follows. First it computes a relation \equiv on $[Q]$ by defining $i \equiv j$ iff $\text{IP}[i, j] \neq \emptyset$, and then makes \equiv an equivalence relation by taking its transitive closure. Then it sorts the partitions of the equivalence relation by their least element, and assign $\hat{\mathbf{x}}[i]$ the index of the partition containing i .

The utility of $\hat{\mathbf{x}}$ is encapsulated in the following lemma, which says that $\hat{\mathbf{x}}$ can work as the “equality pattern” for \mathbf{x} when certain conditions hold.

Lemma 10 *Let $\text{DB} = (\text{ind}_i, W_i)_{i=1}^d$ be a database, $\mathbf{s}, \mathbf{x} \in W^Q$, and let $\hat{\mathbf{x}}$ be the restricted equality pattern. Then for all $i, j \in [Q]$,*

$$\hat{\mathbf{x}}[i] = \hat{\mathbf{x}}[j] \quad \Rightarrow \quad \mathbf{x}[i] = \mathbf{x}[j]$$

and

$$(\mathbf{x}[i] = \mathbf{x}[j]) \wedge (\text{DB}(\mathbf{s}[i]) \cap \text{DB}(\mathbf{s}[j]) \neq \emptyset) \quad \Rightarrow \quad \hat{\mathbf{x}}[i] = \hat{\mathbf{x}}[j].$$

The second assertion is obvious from the definition of $\widehat{\mathbf{x}}$ and IP . For the first assertion, consider i, j such that $\widehat{\mathbf{x}}[i] = \widehat{\mathbf{x}}[j]$. Then there must be indices k_1, \dots, k_t (for some t) with $i = k_1$ and $j = k_t$ and such that the adjacent indices are set equal by the starting definition \equiv . This means that $\text{IP}[k_1, k_2], \text{IP}[k_2, k_3], \dots$ are all non-empty, which by the definition of IP implies that $\mathbf{x}[k_1] = \mathbf{x}[k_2], \mathbf{x}[k_2] = \mathbf{x}[k_3], \dots$, so $\mathbf{x}[i] = \mathbf{x}[k_1] = \mathbf{x}[k_t] = \mathbf{x}[j]$.

The simulator works as follows. On input $(N, \overline{\mathbf{s}}, \text{SP}, \text{RP}, \text{IP}, \mathcal{L}_T(\text{DB}, \mathbf{s}), \mathbf{T}[\mathbf{s}])$, it starts by computing $\widehat{\mathbf{x}}$ as described above. It then proceeds:

```

For  $w \in \widehat{\mathbf{x}}$  and  $\text{ind} \in \bigcup_{i=1} \text{RP}[i]$  do  $H[\text{ind}, w] \xleftarrow{\$} G$ 
For  $w \in \overline{\mathbf{s}}$  do  $\text{WPerms}[w] \xleftarrow{\$} \text{Perm}([\text{SP}[i]])$ 
 $(\text{TSet}, \text{STags}) \leftarrow \mathcal{S}_T(\mathcal{L}_T(\text{DB}, \mathbf{s}), \mathbf{T}[\mathbf{s}])$ 
 $\text{XSet} \leftarrow \emptyset; j \leftarrow 0$ 
For  $w \in \widehat{\mathbf{x}}$  and  $\text{ind} \in \bigcup_{i: \widehat{\mathbf{x}}[i]=w} \text{RP}[i]$  do  $\text{XSet} \leftarrow \text{XSet} \cup \{H[\text{ind}, \widehat{\mathbf{x}}[i]]\}; j \leftarrow j + 1$ 
For  $i = j + 1, \dots, N$  do  $h \xleftarrow{\$} G; \text{XSet} \leftarrow \text{XSet} \cup \{h\}$ .

```

This produces $\text{EDB} = (\text{TSet}, \text{XSet})$. Next it computes the transcript array \mathbf{t} . To fill in $\mathbf{t}[i]$, \mathcal{S} computes the “revealed” indices for that query as $R = \text{RP}[i] \cup \bigcup_{j=1}^Q \text{IP}[i, j]$. It puts the indices in R in canonical order, calling them $(\overline{\text{ind}}_1, \dots, \overline{\text{ind}}_{T'})$. Since every index in R is in $\text{DB}(\mathbf{s}[i])$, we know that $|R| \leq \text{SP}[i]$. The simulator pads $|R|$ up to size $\text{SP}[i]$ by setting $\overline{\text{ind}}_k = \perp$ for $k = T' + 1, \dots, \text{SP}[i]$. It takes the identifiers $(\overline{\text{ind}}_1, \dots, \overline{\text{ind}}_{\text{SP}[i]})$ and computes

```

 $\mathbf{t} \leftarrow \text{TSetRetrieve}(\text{TSet}, \text{stag}[i]); \sigma \leftarrow \text{WPerms}[\overline{\mathbf{s}}[i]]$ 
For  $c = 1, \dots, \text{SP}[i]$  do
   $(y, \mathbf{e}) \leftarrow \mathbf{t}[c]$ 
  If  $\overline{\text{ind}}_{\sigma(c)} \neq \perp$  then  $\text{xtoken}[c] \leftarrow H[\overline{\text{ind}}_{\sigma(c)}, \widehat{\mathbf{x}}[i]]^{1/y}$ 
  Else  $\text{xtoken}[c] \xleftarrow{\$} G$ 
For  $c = \text{SP}[i] + 1, \dots, T$  do  $\text{xtoken}[c] \xleftarrow{\$} G$ .  $\text{Res} \leftarrow \text{ServerSearch}(\text{EDB}, (\text{stag}, \text{xtoken}))$ 
Output  $((\text{stag}, \text{xtoken}), \text{Res}, \text{RP}[i])$ .

```

After computing each $\mathbf{t}[i]$ this way, the simulator outputs $\text{EDB} = (\text{TSet}, \text{XSet})$ and \mathbf{t} .

We claim that the output of \mathcal{S} has the same distribution as the output of Initialize in G_8 . To show this, we first observe that the computed TSet and STags obviously have the same distribution, as they are computed in exactly the same way. To prove the claim it is sufficient show that XSet and all of the xtoken have the same distribution jointly consistent with TSet , STags , and each other.

When computing XSet , \mathcal{S} first initializes $H[\text{ind}, w]$ for each ind that appears as a result and w in $\widehat{\mathbf{x}}$. This differs from G_8 in two ways: First, the w are integers from 1 to Q , not keywords as they are in G_8 , and second, not all ind are initialized. Next it initializes WPerms at locations indexed by entries in $\overline{\mathbf{s}}$, which are again integers instead of keywords.

We claim that both \mathcal{S} and G_8 produce XSet by choosing N uniform and independent group elements. In \mathcal{S} , in the first loop for initializing the XSet , the elements of H indexed by $w \in \widehat{\mathbf{x}}$ and ind that appear as results are added to XSet . It tracks how many additions are done this way with the counter j , and after this loop is add elements until a total of N have been added. In G_8 , one element is added for each $w \in W$ and each ind that contains w (i.e., $\text{ind} \in \text{DB}(w)$), which means that $\sum_{w \in W} |\text{DB}(w)| = N$ elements are added. We stress that in both cases, no entry of H is reused.

We now examine how the xtoken are computed. The values of $\mathbf{t}, y, \mathbf{e}$ are clearly identically distributed (joint with everything so far). We also have that the permutations σ are distributed identically: In both cases they are uniform and independent, except when the same permutation is reused from WPerms . In G_8 this happens when an \mathbf{s} -term is repeated, and it will happen in \mathcal{S} when $\overline{\mathbf{s}}$ repeats, meaning the same pattern of permutations are reused in either case.

Next we observe that in G_8 the identifiers from $\text{DB}(\mathbf{s}[i])$ are used in the random order determined by σ . In the main loop over the identifiers, if that $\sigma(c)$ -th identifier is in either $\text{DB}(\mathbf{s}[i]) \cap \text{DB}(\mathbf{x}[i])$

or in the result set for another query with the same x -term, then the corresponding location in H is read. Otherwise a random group element is chosen. In \mathcal{S} the same logic is used, except the identifiers that do not appear as relevant results are now dummy symbols.

Finally we must show that the entries used from H follow the same “reuse pattern” (when generating the XSet as well as the x token) in G_8 and in \mathcal{S} .

Consider any two index/keyword pairs $(\text{ind}_1, \mathbf{x}[i])$ and $(\text{ind}_2, \mathbf{x}[j])$ read from H by G_8 . By the construction of the game, we know that ind_1 either matches the i -th query or another query with the same x -term, and similarly for ind_2 with the j -th query. The simulator will instead read H at locations $(\text{ind}_1, \hat{\mathbf{x}}[i])$ and $(\text{ind}_2, \hat{\mathbf{x}}[j])$ with the same indices from either RP or IP. We claim that

$$(\text{ind}_1, \mathbf{x}[i]) = (\text{ind}_2, \mathbf{x}[j]) \iff (\text{ind}_1, \hat{\mathbf{x}}[i]) = (\text{ind}_2, \hat{\mathbf{x}}[j]), \quad (1)$$

which will show that the simulation is the same. To prove (1), consider any two such queries. The \Leftarrow direction is easy by Lemma 10. For the \Rightarrow direction, suppose that $(\text{ind}_1, \mathbf{x}[i]) = (\text{ind}_2, \mathbf{x}[j])$. Then we have $\text{ind}_1 = \text{ind}_2$, and also that they are members of the set

$$(\text{RP}[i] \cup \bigcup_{k=1}^Q \text{IP}[i, k]) \cap (\text{RP}[j] \cup \bigcup_{k=1}^Q \text{IP}[j, k])$$

because the games only ever use ind values from these sets as indices. But this means the ind value is in $\text{DB}(\mathbf{s}[i]) \cap \text{DB}(\mathbf{s}[j])$, meaning this intersection is not empty, and thus we have $\hat{\mathbf{x}}[i] = \hat{\mathbf{x}}[j]$ by the second conclusion in Lemma 10.

C Oblivious Cross-Tags Protocol PXT Using Bilinear Pairings

As we saw in the design of OXT in Section 3, the server stores in each tuple corresponding to a keyword w and document xind , a blinded value $y_c = \text{xind} \cdot z_c^{-1}$, where z_c was derived from w and a tuple counter c . This counter served to ensure independent blinding values z_c . Similarly, the Client needed to send an *array* of values $x\text{token}_i$ defined by $x\text{token}[c, i] := g^{F_p(K_X, w_i) \cdot z_c}$.

A natural alternative to using counters to assure independence is to generate $y \leftarrow g^{\text{xind}/z}$, where z is now just derived from w (without any counters). The blinding can be removed from the Client-supplied tokens by using bilinear pairings, if the underlying groups support such pairings. However, we now run into another problem, as the bilinear pairings usually render the DDH assumption invalid, unless we work in bilinear groups $\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_T$, with a pairing $e : \mathcal{G}_1 \times \mathcal{G}_2 \rightarrow \mathcal{G}_T$, and with no easy homomorphism between \mathcal{G}_1 and \mathcal{G}_2 . Such groups do still satisfy DDH, and the assumption is well-known as the SXDH assumption (see [17]).

Thus, the Client can now just send a single $x\text{trap}$ $\mathcal{P}^{F_p(K_X, w_i) \cdot z}$, where \mathcal{P} is a generator of group \mathcal{G}_2 . However, the DDH assumption is not enough to randomize all y values corresponding to a w (i.e. that have the same z in the exponent) as z is also in the exponent of the revealed $x\text{trap}$ for some other w . It would be intriguing to see if one can build an SSE scheme with leakage limited to what OXT leaks by using just the SXDH assumption, although it seems unlikely that it would be as efficient as the one just described.

Surprisingly, we can show that the above scheme, which we call PXT and that is described in detail in Figure 11, is secure with leakage exactly as for OXT by assuming something stronger than SXDH. We call this stronger assumption the augmented SXDH assumption (SAXDH). This augmented assumption remains static, as opposed to many other dynamic assumptions usually employed to achieve efficiency (e.g. the strong Diffie-Hellman assumption [8], q-BDHI assumption etc.). It is also straightforward to verify that it holds in the generic group model (in particular, it is implied by the generalized DH-problem [7]).

EDBSetup(DB)

- Select keys K_S, K_X, K_I, K_Z , and parse DB as $(\text{ind}_i, W_i)_{i=1}^d$.
- Initialize \mathbf{T} to an empty array indexed by keywords from W .
- Initialize XSet to an empty set.
- For each $w \in W$, build the tuple list $\mathbf{T}[w]$ and XSet elements as follows:
 - Initialize \mathbf{t} to be an empty list, and set $K_e \leftarrow F(K_S, w)$.
 - For all ind in $\text{DB}(w)$ in random order, then:
 - * Set $\text{xind} \leftarrow F_p(K_I, \text{ind})$, $z \leftarrow F_p(K_Z, w)$ and $y \leftarrow g^{\text{xind}/z}$.
 - * Compute $\mathbf{e} \leftarrow \text{Enc}(K_e, \text{ind})$, and append (\mathbf{e}, y) to \mathbf{t} .
 - * Set $\text{xtag} \leftarrow e(g, \mathcal{P})^{F_p(K_X, w) \cdot \text{xind}}$ and add xtag to XSet.
 - $\mathbf{T}[w] \leftarrow \mathbf{t}$.
- $(\text{TSet}, K_T) \leftarrow \text{TSetSetup}(\mathbf{T})$.
- Output the key $(K_S, K_X, K_I, K_Z, K_T)$ and $\text{EDB} = (\text{TSet}, \text{XSet})$.

Search protocol

- The client takes an input the key $(K_S, K_X, K_P, K_I, K_Z, K_T)$ and keywords w_1, \dots, w_n to query. It computes $\text{stag} \leftarrow \text{TSetGetTag}(K_T, w_1)$ and then for $i = 2, \dots, n$, $\text{xtoken}_i \leftarrow \mathcal{P}^{F_p(K_Z, w_1) \cdot F_p(K_X, w_i)}$. Sends $(\text{stag}, \text{xtoken}_1, \dots, \text{xtoken}_n)$ to the server.
- The server computes $\mathbf{t} \leftarrow \text{TSetRetrieve}(\text{TSet}, \text{stag})$. For each tuple (\mathbf{e}, y) in \mathbf{t} , if $e(y, \text{xtoken}_i) \in \text{XSet}$ for all $i = 2, \dots, n$, it sends \mathbf{e} to the client.
- The client sets $K_e \leftarrow F(K_S, w_1)$, and for each \mathbf{e} received, the client computes $\text{ind} \leftarrow \text{Dec}(K_e, \mathbf{e})$ and outputs ind .

Figure 11: PXT: BILINEAR PAIRINGS BASED OBLIVIOUS CROSS-TAGS PROTOCOL

C.1 Augmented SXDH Assumption

Consider cyclic groups $\mathcal{G}_1, \mathcal{G}_2$ and \mathcal{G}_T , each of prime order p , with an efficient bilinear pairing $e : \mathcal{G}_1 \times \mathcal{G}_2 \rightarrow \mathcal{G}_T$. The SXDH assumption says that the groups \mathcal{G}_1 and \mathcal{G}_2 are DDH groups (which also implies that the target group \mathcal{G}_T is DDH). For our PXT construction we need a generalization of the SXDH assumption which we call the **symmetric augmented external Diffie-Hellman assumption (SAXDH)**:

Let g and \mathcal{P} be generators of groups \mathcal{G}_1 and \mathcal{G}_2 resp. The **SAXDH** statement requires that the probability that any efficient algorithm can distinguish between the following two distributions is negligible, where both distributions are on $\mathcal{G}_1^4 \times \mathcal{G}_2^4$: the first distribution is obtained by choosing x, s_1, s_2, i and j randomly and independently from \mathbb{Z}_p , and generating the tuple

$$\langle g^{xs_1}, g^{s_1}, g^x, g^{s_2}, \mathcal{P}^{i/s_1}, \mathcal{P}^{i/s_2}, \mathcal{P}^j, \mathcal{P}^{j/s_2} \rangle,$$

and the second distribution is generated identically except that the first element $g^{x_{s1}}$ is replaced by g^y for a random y . The (symmetric) SAXDH assumption also requires that a similar assumption hold when the roles of \mathcal{G}_1 and \mathcal{G}_2 (i.e. g and \mathcal{P}) are reversed.

It is straightforward to see that the SAXDH assumption implies the SXDH assumption.

C.2 PXT Security Theorem

Recall the leakage function \mathcal{L} defined as $\mathcal{L}(\text{DB}, (\mathbf{s}, \mathbf{x})) = (\mathcal{L}_{\text{ext}}(\text{DB}, (\mathbf{s}, \mathbf{x})), \mathcal{L}_T(\mathbf{T}, \mathbf{s}))$ (see the paragraph just before Theorem 9).

Theorem 11 *The 2-conjunctive query SSE scheme PXT over bilinear groups \mathcal{G}_1 , \mathcal{G}_2 and \mathcal{G}_T is \mathcal{L} -semantically-secure against non-adaptive attacks, assuming that the SAXDH assumption holds in \mathcal{G}_1 , \mathcal{G}_2 and \mathcal{G}_T , that F and F_p are secure PRFs, P is a secure PRP, that (Enc, Dec) is an IND-CPA secure symmetric encryption scheme, and that Σ is a (non-adaptively) \mathcal{L}_T -secure and computationally correct T -set implementation.*

The proof for this theorem will be given in the full version of the paper. However, it is instructive to note some key differences from the security proof for OXT. The proof of Theorem 9 (Appendix B) describes how the table H can be populated with $g^{f_X(w) \cdot f_I(\pi(\text{ind}))}$ values, and from this the xtoken can be computed in reverse as $\text{xtoken}[c] \leftarrow H[\text{ind}_{\sigma(c)}, x]^{1/y_c}$.

However, in PXT the values in the XSet are in the target group \mathcal{G}_T , i.e. computed using the pairing, and hence there is no way to go back to get xtoken . Thus, the Simulation is more difficult here, as we describe the simulator in detail next.

C.3 The Simulator for SSE protocol PXT

Let \mathcal{S}_T be the simulator for the \mathcal{L}_T -secure T -set implementation. The SSE ideal game provides the simulator \mathcal{S}_{PXT} with $\mathcal{L}(\text{DB}, (\mathbf{s}, \mathbf{x})) = (\mathcal{L}_{\text{ext}}(\text{DB}, (\mathbf{s}, \mathbf{x})), \mathcal{L}_T(\mathbf{T}, \mathbf{s}))$, which includes the total number of documents d , and sum N of the number of keywords in each document.

The simulator now maintains in addition to table H , two other tables I and Z , and starts by filling all three of them with random and independent values.

It is worth pointing out here that the $\bar{\mathbf{s}}$ does not reveal if two queries have the same x -term. However, from the intersection pattern IP the following matrix M_x of size $Q \times Q$ can be obtained: $M_x[i, j] = 1$ iff $S_i \cap S_j \neq \emptyset$. Now note, $M_x[i, j] = 1$ iff $\mathbf{q}[i]_x = \mathbf{q}[j]_x$ and $I_\pi(\mathbf{s}[i]) \cap I_\pi(\mathbf{s}[j]) \neq \emptyset$. Thus, M_x reveals partial repetition pattern of x -terms conditional on IP. Similar to profile $\bar{\mathbf{s}}$, one can also obtain a vector $\bar{\mathbf{x}}$ from M_x as follows:

- Define $\mu_x(i) = j$, where j is the least query number such that $M_x[i, j] = 1$.
- Define $\bar{\mathbf{x}}[i] = z$, where there are exactly z queries j , $j \leq i$, s.t. $\mu_x(j) = j$.

Using the leakage profiles IP and RP, \mathcal{S}_{PXT} does the following.

1. For each j in the range of $\bar{\mathbf{s}}$, define $CCC_j = \bigcup_{i: \bar{\mathbf{s}}[i]=j} (R_i \cup S_i)$.
2. For each $i \in [1..Q]$, it generates set $\mathbf{t}'[i] = \{g^{I[r]/Z[\bar{\mathbf{s}}[i]]}\}_{r \in CCC_{\bar{\mathbf{s}}[i]}}$.
3. It generates $\mathbf{t}[i]$ to be $\mathbf{t}'[i]$ union a set of random and independent values such that the total number of entries in $\mathbf{t}[i]$ is $\text{SP}[i]$ (which is same for all j s.t. $\bar{\mathbf{s}}[i] = \bar{\mathbf{s}}[j]$).

4. The simulator \mathcal{S}_{PXT} calls \mathcal{S}_T with $\mathcal{L}_T(\text{DB}, \mathbf{s})$ and $\langle \mathbf{t}[i] \rangle_{i \in [1..Q]}$, which in turn returns TSet and STags , which \mathcal{S}_{PXT} outputs.
5. As a simulation of the i -th xtrap, \mathcal{S}_{PXT} outputs $\text{xtoken}[i] = \mathcal{P}^{Z[\overline{\mathbf{S}}[i]] \cdot H[\overline{\mathbf{x}}[i]]}$.
6. As a simulation of the XSet , \mathcal{S}_{PXT} outputs the indicator function of the union of the following sets: for each $i \in [1..Q]$, the set consists of $e(g^{I[r]/Z[\overline{\mathbf{S}}[i]]}, \text{xtoken}[i])$ for each r in R_i .

In the full version of the paper we will show that under the SAXDH assumption, the view of the adversary in the SSE ideal game using \mathcal{S}_{PXT} is indistinguishable from the view of the adversary in the real game using PXT .