

TWO-IN-one-SSE: Fast, Scalable and Storage-Efficient Searchable Symmetric Encryption for Conjunctive and Disjunctive Boolean Queries

Arnab Bag
Indian Institute of Technology
Kharagpur, West Bengal, India
arnabbag@iitkgp.ac.in

Debadrita Talapatra
Indian Institute of Technology
Kharagpur, West Bengal, India
debadritat.fg2219@kgpian.iitkgp.ac.in

Ayushi Rastogi
Indian Institute of Technology
Kharagpur, West Bengal, India
rayushi835@kgpian.iitkgp.ac.in

Sikhar Patranabis*
IBM Research India
Bengaluru, Karnataka, India
sikhar.patranabis@ibm.com

Debdeep Mukhopadhyay
Indian Institute of Technology
Kharagpur, West Bengal, India
debdeep@cse.iitkgp.ac.in

ABSTRACT

Searchable Symmetric Encryption (SSE) supports efficient yet secure query processing over outsourced symmetrically encrypted databases without the need for decryption. A longstanding open question has been the following: can we design a fast, scalable, linear storage and low-leakage SSE scheme that efficiently supports arbitrary Boolean queries over encrypted databases? In this paper, we present the design, analysis and prototype implementation of the first SSE scheme that efficiently supports conjunctive, disjunctive *and* more general Boolean queries (in both the conjunctive and disjunctive normal forms) while scaling smoothly to extremely large encrypted databases, and while incurring *linear* storage overheads and supporting extremely fast query processing in practice. We quantify the leakage of our proposal via a rigorous cryptographic analysis and argue that it achieves security against a well-known class of leakage-abuse and volume analysis attacks. Finally, we demonstrate the storage-efficiency and scalability of our proposed scheme by presenting experimental results of a prototype implementation of our scheme over large real-world databases.

KEYWORDS

Searchable Encryption, Database Encryption, Encrypted Search, Boolean Queries, Conjunctive and Disjunctive Queries.

1 INTRODUCTION

The advent of cloud computing potentially allows individuals and organizations to outsource storage and processing of large volumes of data to third party servers. However, this leads to concerns surrounding the confidentiality of the data stored on the third party cloud servers and outsourced access for processing.

Consider, for instance, a client that offloads an encrypted database of (potentially sensitive) emails to an untrusted cloud server.

*Also with ETH Zürich, Switzerland and Visa Research USA (part of the work was done while the author was affiliated with these institutions).

This work is licensed under the Creative Commons Attribution 4.0 International License. To view a copy of this license visit <https://creativecommons.org/licenses/by/4.0/> or send a letter to Creative Commons, PO Box 1866, Mountain View, CA 94042, USA.

Proceedings on Privacy Enhancing Technologies 2023(1), 115–139

© 2023 Copyright held by the owner/author(s).

<https://doi.org/10.56553/popets-2023-0008>



At a later point of time, the client might want to issue a query of the form: *retrieve all emails received from xyz@foobar.org OR abc@foobar.org AND with "research" in the subject field*. Ideally, the client should be able to perform this task without revealing any sensitive information to the server, such as the sources and contents of the emails, the keywords underlying a given query, the distribution of keywords across emails, etc. Unfortunately, existing techniques such as Fully Homomorphic Encryption (FHE) [14] and Oblivious RAM (ORAM) [16], that potentially support such an “ideal” notion of privacy, are currently unsuitable for wide-scale practical deployment due to high performance and storage overheads.

Searchable Symmetric Encryption. Searchable Symmetric Encryption (SSE) [3, 4, 6–13, 15, 19, 21–23, 29, 30] is the study of provisioning symmetric-key encryption schemes with search capabilities. The goal of SSE is two-fold: (a) to allow a (potentially untrusted) server to execute keyword search queries directly on a collection of a client’s encrypted documents in an efficient manner, and (b) to ensure client privacy by minimising the amount of information “leakage” to the server in the process. Some examples of leakage include the database size, *query pattern* (which queries correspond to the same keyword) and the *access pattern* (the set of file identifiers matching a given query).

SSE for Boolean Queries. The example query over an email database that we outlined above is an instance of what we call a *Boolean query*, in the sense that it can be viewed as a Boolean formula involving certain equality predicates over keywords, connected by AND and OR operators. In this paper, we broadly investigate the following question:

Can we design a fast, scalable, storage-efficient and low-leakage SSE scheme for general Boolean queries?

This seemingly natural question has, somewhat surprisingly, also been a longstanding open question. In particular, while significant progress has been made in designing efficient SSE schemes for simpler sub-classes of Boolean queries (such as atomic equality predicates and conjunctions of keywords), the handful of existing SSE schemes supporting disjunctive and general Boolean queries incur extremely large encrypted storage overheads (quadratic in the size of the plain database), which makes them impractical for real-world deployment. We briefly summarise the state-of-the-art

on both conjunctive and disjunctive SSE below followed by main contributions of this work.

1.1 Background and Related Work

Initial constructions of SSE focused on single keyword search and their related extensions. Recent SSE algorithms now support multi-keyword search. We briefly go over the current status of multi-keyword SSE constructions below.

SSE for Conjunctions. In a seminal work [7], Cash et al. proposed Oblivious Cross Tags (OXT) - an efficient, highly scalable and low-leakage (not leaking more than benign information) SSE schemes supporting conjunctive keyword queries over encrypted document collections. Since then, a number of SSE schemes supporting conjunctive keyword queries with a variety of leakage vs efficiency trade-offs have been proposed in different settings [6, 23, 29]. Unfortunately, these schemes are neither efficient nor low-leakage when processing disjunctive or general Boolean queries. For example, when processing a disjunctive query over (w_1, \dots, w_n) , these solutions leak the set of documents matching each w_i *individually*, which can be devastating in the face of existing leakage-abuse attacks [2, 5, 25, 32].

SSE for Disjunctions. While SSE for single and conjunctive keyword queries has been studied quite extensively, SSE for disjunctive queries has received much less attention. To the best of our knowledge, only IEX-2LEV and IEX-ZMF due to Kamara and Moataz [19] support reasonably efficient query processing without incurring potentially devastating leakage.

The IEX family of schemes [19, 26] has a few disadvantages that we outline here. First, its performance for conjunctive queries over real-world databases is significantly worse as compared to OXT. Secondly, it is incompatible with OXT and its follow-up schemes [6, 23, 29]; so it does not lead to a common solution that supports both conjunctive and disjunctive queries efficiently. Finally, the IEX family of schemes incurs a (worst-case) storage overhead that grows *quadratically* with the number of keywords in the database. This makes it impractical for deployment over real-world databases.

1.2 Our Contributions

In this paper, we present the design, analysis and prototype implementation of the first SSE scheme that efficiently supports conjunctive, disjunctive *and* more general (and complex) Boolean queries (in both the conjunctive and disjunctive normal forms) while scaling smoothly to extremely large encrypted databases, and while incurring *linear* storage overheads and little query processing overheads in practice. Our scheme is named TWo-IN-one-SSE, or TWINSSE in short. We expand on our contributions and techniques below.

Supporting Conjunctive “and” Disjunctive Queries. Our core technical contribution is a novel mechanism for designing SSE schemes that support both conjunctive and disjunctive keyword searches in a fully compatible manner. At a high level, we achieve this as follows. Given *any* conjunctive SSE scheme (i.e., any generic SSE scheme that only supports conjunctive queries), we present a generic black-box transformation that yields an SSE scheme supporting conjunctive, disjunctive and general Boolean queries in the conjunctive normal form (CNF) and disjunctive normal form (DNF). Our transformation does not rely on any special properties of the

underlying conjunctive SSE scheme. This allows it to be instantiated from *any* existing conjunctive SSE scheme, including OXT [7]. To the best of our knowledge, such a generic transformation from a conjunctive SSE scheme to an SSE scheme for general and complex Boolean queries has not been studied before in the SSE literature¹.

A Naïve Approach. The naïve approach for supporting disjunctive queries in a generic way using a system that only supports conjunctive queries is to allow for “negative searches”, wherein given a keyword, we can efficiently retrieve the set of documents that *do not* contain the keyword. Then, given a disjunctive query of the form $q = (w_1 \vee \dots \vee w_n)$, we can transform q into the following conjunctive query $\bar{q} = (\bar{w}_1 \wedge \dots \wedge \bar{w}_n)$, where for each i , \bar{w}_i denotes the “negated keyword” that (hypothetically) occurs in every document that does not contain w_i ; consequently, \bar{q} can be viewed as the negated counterpart to the original query q .

This approach has two major disadvantages. First, it requires us to design data structures that support efficiently retrieving, for each keyword, not only the set of documents it occurs in, but also the set of documents that it does not occur in, while maintaining data and query privacy. This is likely to lead to massive blowup in storage. Secondly, and most crucially, for disjunctive queries involving less frequent keywords (which is what we expect from a very large proportion of the queries), the overall computational and communication complexity suffers a huge blowup, since it is now proportional to the result set for the negated query, which would be almost the entire set of documents. Our aim is to design a generic transformation mechanism that is significantly more efficient. To this end, we introduce and use a novel concept called *meta-keywords*.

Using Meta-Keywords. The technical centrepiece of our generic transformation is the concept of *meta-keywords*, which we introduce in this paper. At a high level, a meta-keyword mkw_i is a disjunction of certain carefully chosen keywords of the form $mkw_i = (w_{i_1} \vee w_{i_2} \vee \dots \vee w_{i_r})$, that we pre-process and store at setup in an inverted search index².

Our core technical observation is the following: given a database with N keywords, there exists an $O(N)$ -sized set \mathcal{S} of meta-keywords such that for *any* disjunctive query of the form $q = (w_1 \vee \dots \vee w_n)$ for $n \leq N$, there exists a *meta-query* which is a *conjunction over meta-keywords* of the form $q' = (mkw_1 \wedge \dots \wedge mkw_n)$ such that $mkw_1, \dots, mkw_n \in \mathcal{S}$, and such that $\text{Result-Set}(q) \subseteq \text{Result-Set}(q')$.

The non-triviality of our approach lies in addressing the following challenges simultaneously:

- **Coverage:** Designing an $O(N)$ -sized meta-keyword set that “covers” an $O(2^N)$ -sized space of all possible disjunctive queries.
- **Efficiency:** Minimizing the overheads due to filtering of “spurious” documents in the result “meta-set” (i.e. ensuring that $\text{Result-Set}(q')$ is as close to $\text{Result-Set}(q)$ as possible).

¹We note here that OXT does support Boolean queries beyond simple conjunctions, albeit where the query must be in a restricted *searchable normal form* (SNF) [7]; our transformation is significantly more general in the sense that it extends to any CNF or DNF formula over keywords, well beyond the scope of SNF queries.

²An inverted search index is a data structure popularly used by essentially all SSE schemes that is indexed by the keywords, and stores, for each keyword, the set of documents it contains, albeit in encrypted form.

- **Security:** Minimizing leakage by ensuring that the meta-keywords reveal as little information as possible about the underlying keywords being actually queried.

While achieving these requirements simultaneously appears challenging at first sight, we develop a systematic and formal approach that allows us to achieve them for *any* database; more formally, given a database, we show how to convert the same into a *meta-database* equipped with a linearly-sized set of meta-keywords that meets all of the aforementioned requirements. We formalise these properties in Sections 3 and 4.

TWINSSE. We use the aforementioned transformation to design our overall solution, that we call TWO-IN-one-SSE, or TWINSSE in short. As the astute reader might have already guessed, given a conjunctive SSE scheme, our design of TWINSSE uses the following two-step approach:

- **Step-1:** Given any database **DB**, convert it into the corresponding meta-database $\widehat{\mathbf{DB}}$, where $\widehat{\mathbf{DB}}$ can be viewed as a database equipped with two kinds of keywords – the original keywords and the meta-keywords.
- **Step-2:** Apply the conjunctive SSE scheme to encrypt and query the database $\widehat{\mathbf{DB}}$.

Note that conjunctive query processing over $\widehat{\mathbf{DB}}$ proceeds exactly as it would over **DB**, and requires no additional query planning on the part of the client. Disjunctive query processing is more involved because it requires the client to plan the meta-query. In Section 4, we formally describe how this can be done using $O(n)$ computation (which is the information-theoretic minimum for any n -word disjunctive query). Finally, in Appendix G, we describe a hybrid query planning approach that allows handling general Boolean queries in CNF and DNF expressions in an efficient manner. We note here that the ability to efficiently handle CNF and DNF queries effectively allows TWINSSE to handle complex Boolean queries (involving both conjunctive and disjunctive clauses) by casting them into either CNF or DNF formulae over keywords.

We then present a concrete instantiation of TWINSSE from OXT as the baseline conjunctive SSE scheme. We denote this version of OXT as TWINSSE_{OXT}. Details of TWINSSE_{OXT} and, in particular, its handling of disjunctive queries, are presented in Section 4.3. Additionally, we present an elaborate discussion on executing complex Boolean queries using TWINSSE_{OXT} in Appendix G.

Leakage Analysis. We formally detail the leakage profile of TWINSSE_{OXT} in Appendix D. In order to analyse the impact of this leakage on the security of TWINSSE_{OXT}, we perform a detailed cryptanalysis of TWINSSE_{OXT} in Appendix F. In particular, we show that known leakage-abuse attacks [5, 32], volume analysis-based attacks [2], and the state-of-the-art SAP attack [25] fail against TWINSSE_{OXT} in practical adversarial settings.

Experimental Evaluation. We present a C++ implementation of TWINSSE_{OXT} along with performance figures in Section 5. We experimented over the Enron email corpus^{3,4} for compatibility with previous SSE literature. The data set contains around 170K keywords, 500K documents and 20 million unique keyword-document

pairs. Our experiments validate that TWINSSE_{OXT} supports extremely fast conjunctive, disjunctive and more complex Boolean queries (in CNF/DNF expression), and substantially outperforms the IEX family of schemes on two counts: (a) storage requirements for the encrypted database, and (b) practical search performance for conjunctive queries, while also achieving comparable practical search performance for disjunctive queries.

2 PRELIMINARIES AND BACKGROUND

In this section we introduce the notations used in the rest of the paper, as well as preliminary background material on SSE.

2.1 Notations

We write $x \stackrel{R}{\leftarrow} \mathcal{X}$ to represent that an element x is sampled uniformly at random from a set/distribution \mathcal{X} . The output x of a deterministic algorithm \mathcal{A} is denoted by $x = \mathcal{A}$ and the output x' of a randomized algorithm \mathcal{A}' is denoted by $x' \leftarrow \mathcal{A}'$. For $a, b \in \mathbb{Z}$ such that $a, b \geq 0$, we denote by $[a]$ and $[a, b]$ the set of integers lying between 1 and a (both inclusive), and the set of integers lying between a and b (both inclusive), respectively.

Databases. Let $\Delta = \{w_1, \dots, w_N\}$ be a dictionary of keywords, and let $\mathcal{F} = \{f_1, \dots, f_D\}$ be a collection of documents, such that each document f_i is associated with a unique identifier id_i and contains keywords from Δ . We assume that standard set operations including union and intersection are allowed over Δ . We denote by **DB** a database of identifier-keyword pairs, such that $(id, w) \in \mathbf{DB}$ if and only if the document with identifier id contains the keyword w . We denote by $\mathbf{DB}(w)$ the set of all identifiers corresponding to documents containing w . We denote by $|\Delta|$ the number of distinct keywords in **DB**, by $|\mathbf{DB}|$ the number of distinct $id - w$ pairs in **DB**, and by $|\mathbf{DB}(w)|$ the number of documents containing w .

Conjunctive and Disjunctive Queries. We represent a conjunctive query over n distinct keywords w_1, \dots, w_n as $q = w_1 \wedge \dots \wedge w_n$ and define the set $\mathbf{DB}(q)$ as $\mathbf{DB}(q) = \bigcap_{i=1}^n \mathbf{DB}(w_i)$. Similarly, we represent a disjunctive query over n distinct keywords w_1, \dots, w_n as $q = (w_1 \vee \dots \vee w_n)$ and define the set $\mathbf{DB}(q)$ as $\mathbf{DB}(q) = \bigcup_{i=1}^n \mathbf{DB}(w_i)$.

Throughout the paper, we use the notation $\mathcal{R}_q = \mathbf{DB}(q)$ to represent the result of searching a query q (irrespective of the query type), unless otherwise specified.

2.2 Searchable Symmetric Encryption

Any SSE scheme [7, 11] consists of a polynomial-time algorithm **SETUP** executed by the client, and an interactive protocol **SEARCH** executed jointly by the client and the server:

- **SETUP**($1^\lambda, \mathbf{DB}$): Takes as input the security parameter λ and a database **DB**, and outputs the tuple (sk, st, \mathbf{EDB}) , where sk is the client’s secret-key, st is the client’s internal state, and **EDB** is the encrypted database.
- **SEARCH**(sk, st, q, \mathbf{EDB}): The client takes as input the secret-key sk , its state st and a query q , while the server takes as

³<https://www.cs.cmu.edu/~enron/>

⁴<https://www.kaggle.com/wcukierski/enron-email-dataset>

input the encrypted database **EDB**. At the end of the protocol, the client outputs $\mathbf{DB}(q)^5$.

Correctness. An SSE scheme is said to be correct if for every database **DB** and for every query q , the output of the SEARCH protocol contains $\mathbf{DB}(q)$ with overwhelming probability.

Security. We refer to [7, 11] for the standard simulation security definition of SSE against semi-honest adversaries in the real world-ideal world paradigm.

3 TWIN SSE: SIMPLIFIED VERSION

In this section, we introduce TWIN SSE. For ease of representation, we first present a simplified version, which we refer to as TWIN SSE_{BASIC}.

3.1 The Core Tool: Meta-Keywords

We begin by describing the core technical tool for our construction, which we refer to as *meta-keywords*. Let $\Delta = \{w_1, \dots, w_N\}$ be a dictionary of keywords (N is total number of keywords in **DB**), and assume (without loss of generality) that these keywords are arranged in increasing order of frequency, i.e.,

$$|\mathbf{DB}(w_1)| \leq |\mathbf{DB}(w_2)| \leq \dots \leq |\mathbf{DB}(w_N)|.$$

Super-Keyword. We begin by defining a super-keyword, which is simply a *disjunction* of some subset of the keywords in Δ . Formally, a super-keyword \bar{w} is represented by a bit-string of the form $\bar{w} = (b_1, b_2, \dots, b_N) \in \{0, 1\}^N$, such that

$$\mathbf{DB}(\bar{w}) = \bigcup_{\ell \in [N] \text{ s.t. } b_\ell = 1} \mathbf{DB}(w_\ell).$$

Note that one could equivalently represent \bar{w} using the actual constituent keywords; we use the bit-string representation because it makes the description of our strategy easier to follow, and also more efficiently implementable.

Meta-Keyword. We now define a meta-keyword. At a high level, a meta-keyword is a “special” super-keyword with a *single contiguous stretch of 0-entries* in its Boolean representation. Formally, a meta-keyword is defined as follows.

DEFINITION 3.1 (META-KEYWORD). A meta-keyword $\text{mkw}_{i,j}$ is a super-keyword indexed by $(i, j) \in [N] \times [N]$ such that $i \leq j$, represented as bit-string $\text{mkw}_{i,j} = (b_1, b_2, \dots, b_N) \in \{0, 1\}^N$, where for each $\ell \in [N]$, we have

$$b_\ell = \begin{cases} 0 & \text{if } \ell \in [i, j], \\ 1 & \text{otherwise.} \end{cases}$$

In other words, for the meta-keyword $\text{mkw}_{i,j}$, we have

$$\mathbf{DB}(\text{mkw}_{i,j}) = \bigcup_{\ell \in [N] \setminus [i,j]} \mathbf{DB}(w_\ell).$$

Informally, one can view a meta-keyword $\text{mkw}_{i,j}$ as a disjunction over $\Delta = \{w_1, \dots, w_N\}$ *excluding* a contiguous sequence of keywords $(w_i, w_{i+1}, \dots, w_j)$.

We also let $\text{mkw}^* = 1^N$ denote the special “all-ones” meta-keyword. Finally, let $\mathcal{S}_{\text{mkw}, \Delta} = \{\text{mkw}_{i,j}\}_{i \leq j} \cup \{\text{mkw}^*\}$ be the set

⁵We also make the implicit assumption that upon obtaining the set of document identifiers corresponding to a query, the client performs an additional interaction with the server to actually retrieve the documents with these identifiers.

of all meta-keywords over the dictionary Δ . It is easy to see that for $|\Delta| = N$, we have $|\mathcal{S}_{\text{mkw}, \Delta}| = O(N^2)$.

Using a Meta-Keyword. A reader might wonder why we choose the above definition of a meta-keyword. To begin with, note that while pre-processing and storing an inverted search index consisting of all super-keywords along with the original keywords allows us to trivially answer conjunctive, disjunctive, and more general Boolean queries in a fully compatible manner. However, this would require exponential storage, and therefore is not practically feasible.

Hence, our approach is to look for a poly-sized subset of the set of all possible super-keywords that, if pre-processed and stored as part of the inverted search index, would allow us to “cover” any disjunctive query. It turns that the set of meta-keywords is indeed this set. We make this explicit by stating the following (informal) claim. We subsequently make this claim more formal and prove it.

CLAIM 3.1 (INFORMAL). For any disjunctive query of the form $q = (w_{\ell_1} \vee \dots \vee w_{\ell_n})$ (where $n \leq N$), there exists a meta-query which is a conjunction over meta-keywords of the form

$$q_{\text{mkw}} = (\text{mkw}_{i_1, j_1} \wedge \dots \wedge \text{mkw}_{i_n, j_n}),$$

such that $\text{mkw}_{i_1, j_1}, \dots, \text{mkw}_{i_n, j_n} \in \mathcal{S}_{\text{mkw}, \Delta}$, and such that

$$\begin{aligned} \mathbf{DB}(q) &= \bigcup_{k \in [n]} \mathbf{DB}(w_{\ell_k}) \subseteq \bigcap_{k \in [n]} \mathbf{DB}(\text{mkw}_{i_k, j_k}) \\ &= \mathbf{DB}(q_{\text{mkw}}). \end{aligned}$$

Using Claim 3.1 (Overview). As an astute reader might have already observed, this claim allows us to convert *any* disjunctive query over the original set of keywords into a conjunctive query over the set of meta-keywords. Consequently, given a database **DB** over a dictionary Δ , suppose we pre-process **DB** at setup to construct the set of meta-keywords $\mathcal{S}_{\text{mkw}, \Delta}$, and build an augmented meta-database $\widehat{\mathbf{DB}}$ over the meta-dictionary $\widehat{\Delta} = (\Delta \cup \mathcal{S}_{\text{mkw}, \Delta})$ (consisting of *both* the original keywords *and* the meta-keywords). We can then use this augmented (plaintext) database together with *any* conjunctive SSE scheme in a black-box manner to build an SSE scheme that supports both conjunctive *and* disjunctive queries. The only price we pay is the $O(N^2)$ storage overhead; we subsequently show how to reduce this to $O(N)$.

3.2 Meta-Keywords as “Covering” Set

We now formalize and prove Claim 3.1. Before detailing the formal proof, we illustrate why this claim is true via a simple toy-example.

Toy-Example. Consider a database **DB** with 10 documents (indexed as $\{\text{id}_1, \dots, \text{id}_{10}\}$) over a four-keyword dictionary

$$\Delta = \{w_1, w_2, w_3, w_4\},$$

such that

$$\begin{aligned} \mathbf{DB}(w_1) &= \{\text{id}_5, \text{id}_{10}\}, \\ \mathbf{DB}(w_2) &= \{\text{id}_2, \text{id}_5, \text{id}_8, \text{id}_9\}, \\ \mathbf{DB}(w_3) &= \{\text{id}_1, \text{id}_2, \text{id}_4, \text{id}_7, \text{id}_9\}, \\ \mathbf{DB}(w_4) &= \{\text{id}_1, \text{id}_2, \text{id}_3, \text{id}_4, \text{id}_5, \text{id}_6, \text{id}_8, \text{id}_{10}\}. \end{aligned}$$

Now consider the following example disjunctive queries q and the corresponding meta-queries q_{mkw} .

Example-1. Let $q = (w_1 \vee w_2 \vee w_3)$, and $q_{\text{mkw}} = \text{mkw}_{4,4}$, where $\text{mkw}_{4,4} = (1, 1, 1, 0)$.

$$\begin{aligned} \text{DB}(\text{mkw}_{4,4}) &= \bigcup_{\ell \neq 4} \text{DB}(w_\ell) \\ &= \{\text{id}_1, \text{id}_2, \text{id}_4, \text{id}_5, \text{id}_7, \dots, \text{id}_{10}\}, \\ \text{DB}(q_{\text{mkw}}) &= \text{DB}(\text{mkw}_{4,4}). \end{aligned}$$

Now, we have

$$\text{DB}(q) = \{\text{id}_1, \text{id}_2, \text{id}_4, \text{id}_5, \text{id}_7, \dots, \text{id}_{10}\} = \text{DB}(q_{\text{mkw}}).$$

Example-2. Alternatively, suppose that $q = (w_2 \vee w_3)$ and $q_{\text{mkw}} = (\text{mkw}_{1,1} \wedge \text{mkw}_{4,4})$, where $\text{mkw}_{1,1} = (0, 1, 1, 1)$ and $\text{mkw}_{4,4} = (1, 1, 1, 0)$.

$$\begin{aligned} \text{DB}(\text{mkw}_{1,1}) &= \bigcup_{\ell \neq 1} \text{DB}(w_\ell) = \{\text{id}_1, \dots, \text{id}_{10}\}, \\ \text{DB}(q_{\text{mkw}}) &= \text{DB}(\text{mkw}_{1,1}) \cap \text{DB}(\text{mkw}_{4,4}) \\ &= \{\text{id}_1, \text{id}_2, \text{id}_4, \text{id}_5, \text{id}_7, \dots, \text{id}_{10}\}. \end{aligned}$$

Now, we have

$$\text{DB}(q) = \{\text{id}_1, \text{id}_2, \text{id}_4, \text{id}_5, \text{id}_7, \text{id}_8, \text{id}_9\} \subseteq \text{DB}(q_{\text{mkw}}).$$

Formal Statement. We now state the following formal version of Claim 3.1.

LEMMA 3.1. *Let $q = (w_{\ell_1} \vee \dots \vee w_{\ell_n})$ for some $n \leq N$, where $\ell_1 \leq \dots \leq \ell_n$, and let*

$$q_{\text{mkw}} = (\text{mkw}_{i_0, j_0} \wedge \text{mkw}_{i_1, j_1} \wedge \dots \wedge \text{mkw}_{i_n, j_n}),$$

where for each $k \in [0, n-1]$, we have

$$\text{mkw}_{i_k, j_k} = \begin{cases} \text{mkw}_{\ell_{k+1}, \ell_{k+1}-1} & \text{if } \ell_{k+1} > \ell_k + 1, \\ \phi & \text{otherwise,} \end{cases}$$

where we define $\ell_0 := 0$ and ϕ denotes an “empty” meta-keyword, and

$$\text{mkw}_{i_n, j_n} = \begin{cases} \text{mkw}_{\ell_n+1, n} & \text{if } \ell_n < n, \\ \text{mkw}^* & \text{otherwise.} \end{cases}$$

Then we have

$$\text{DB}(q) = \bigcup_{k \in [n]} \text{DB}(w_{\ell_k}) \subseteq \text{DB}(q_{\text{mkw}}).$$

Observe that for the specific examples stated above, the conjunctive meta-query q_{mkw} exactly follows the generic conjunctive meta-query laid out in the above lemma. Here, ℓ_k denotes the index of k 'th query keyword in Δ , whereas (i_k, j_k) denote the start and end indices of the absent keywords stretch in each mkw . Following the above mkw formulation, we see that mkw^* occurs only if all keywords in Δ are present in the query q - when the value of n (number of keywords in q) is same as the number of keywords in Δ . In this case, the index of the last query keyword in Δ (or ℓ_k) is equal to the number of query keywords n , and only mkw^* is selected. This is a rather unusual case that rarely occurs in real applications.

Note that, at a high level, to prove this lemma it suffices that to prove that for each keyword w_{ℓ_k} in the queried disjunction q , we have $\text{DB}(w_{\ell_k}) \subseteq \text{DB}(q)$. In more detail, we would like to prove that for each keyword w_{ℓ_k} in the queried disjunction q , we have $\text{DB}(w_{\ell_k}) \subseteq \text{DB}(\text{mkw}_{i_{\bar{k}}, j_{\bar{k}}})$ for each $\bar{k} \in [0, n]$. The proof follows from the fact that for each $k \in [n]$, the following must be true:

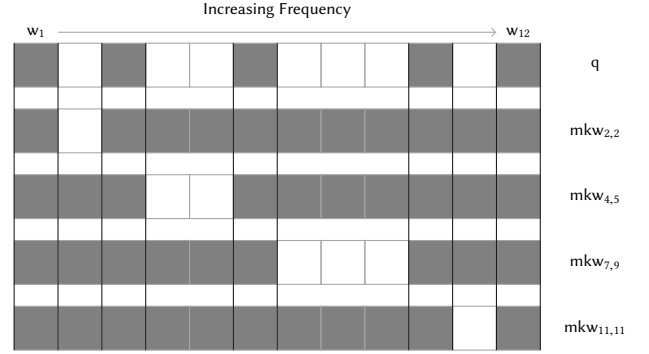


Figure 1: Expressing disjunctive query q in terms of mkw -s with a single stretch of 0s. In this example, $\Delta = \{w_1, \dots, w_{12}\}$, and $q = w_1 \vee w_3 \vee w_6 \vee w_{10} \vee w_{12}$ (where $\ell_1 = 1, \ell_2 = 3, \ell_3 = 6, \ell_4 = 10, \ell_5 = 12$ and $n = 5$). Note that, each mkw has w s present at the same places where a w is present in q . The stretches of 0s (absence of w s) ensure that when the mkw s are ANDed together (searched in a conjunctive manner), only w s in the original query q remain. (Gray and white cells represent 1 and 0, respectively.)

- The index corresponding to the keyword w_{ℓ_k} has a 1-entry in every non-empty meta-keyword in the set $\{\text{mkw}_{i_{\bar{k}}, j_{\bar{k}}}\}_{\bar{k} < k}$. This is because the “stretch” of 0-entries in each such meta-keyword ends *before* the index ℓ_k .
- The index corresponding to the keyword w_{ℓ_k} has a 1-entry in every non-empty meta-keyword in the set $\{\text{mkw}_{i_{\bar{k}}, j_{\bar{k}}}\}_{\bar{k} \geq k}$. This is because the “stretch” of 0-entries in each such meta-keyword starts *after* the index ℓ_k .
- Finally, the index corresponding to the keyword w_{ℓ_k} has, by default, a 1-entry in mkw^* - the “all-ones” meta-keyword.

Combining these observations, we get that for each keyword w_{ℓ_k} in the queried disjunction q , we have $\text{DB}(w_{\ell_k}) \subseteq \text{DB}(\text{mkw}_{i_{\bar{k}}, j_{\bar{k}}})$ for each $\bar{k} \in [0, n]$, as desired. Figure 1 captures the aforementioned intuition pictorially. We defer the formal proof of Lemma 3.1 to Appendix A due to space constraints.

3.3 TWINSSE_{BASIC}

We now put everything together in our basic scheme TWINSSE_{BASIC}. Let CSSE = (CSSE.SETUP, CSSE.SEARCH) be *any* generic conjunctive SSE scheme. Given CSSE, we construct

$$\text{TWINSSE}_{\text{BASIC}} = \begin{cases} \text{TWINSSE}_{\text{BASIC}}.\text{SETUP} \\ \text{TWINSSE}_{\text{BASIC}}.\text{SEARCH} \end{cases}$$

as described subsequently. Our description here is slightly informal due to space constraint, but captures the overall idea of our approach. More details are available with the final construction in Section 4. We also present brief details of processing purely conjunctive or purely disjunctive; we defer the discussion on our treatment of general Boolean formulae to Appendix G.

TWINSSE_{BASIC}.SETUP($1^\lambda, \text{DB}$): Given a database DB over a dictionary Δ , construct the set of meta-keywords $\mathcal{S}_{\text{mkw}, \Delta}$ as described

above. Let $\widehat{\mathbf{DB}}$ denote the meta-database over $\widehat{\Delta} = (\Delta \cup \mathcal{S}_{\text{mkw},\Delta})$ (consisting of *both* the original keywords *and* the meta-keywords). Output

$$(\text{sk}, \text{st}, \widehat{\mathbf{EDB}}) \leftarrow \text{CSSE.SETUP}(1^\lambda, \widehat{\mathbf{DB}}).$$

$\text{TWINSSE}_{\text{BASIC}}.\text{SEARCH}(\text{sk}, \text{st}, q; \widehat{\mathbf{EDB}})$: Given a query q , proceed as follows:

- If q is a purely conjunctive query, output

$$\mathbf{DB}(q) = \text{CSSE.SEARCH}(\text{sk}, \text{st}, q; \widehat{\mathbf{EDB}}).$$

- If q is a purely disjunctive query, construct the conjunctive meta-query q_{mkw} as described in Lemma 3.1, which allows the client to recover

$$\mathbf{DB}(q_{\text{mkw}}) = \text{CSSE.SEARCH}(\text{sk}, \text{st}, q; \widehat{\mathbf{EDB}}),$$

and locally filter $\mathbf{DB}(q) \subseteq \mathbf{DB}(q_{\text{mkw}})$.

Correctness. Correctness is immediate from Lemma 3.1 (see Section A) and correctness of the CSSE scheme. However, this includes the trivial case of returning entire database upon searching a disjunctive query. To avoid such trivial inclusions, we bound the returned result set size close to the actual result set via a *precision* parameter. We define precision η by the following ratio.

$$\eta = \frac{|\mathbf{DB}(q)|}{|\mathbf{DB}(q_{\text{mkw}})|}$$

At a high level, this precision parameter η is a measure of the fraction of spurious ids present in the obtained result set compared to the actual result set. Thus, the correctness can now be defined by the following statement.

For a functionally correct and exact⁶ conjunctive SSE scheme CSSE, a plaintext database \mathbf{DB} , and a disjunctive query q with the corresponding transformed meta-query q_{mkw} , TWINSSE is functionally correct if the following expressions hold.

$$\text{sk}, \text{st}; \widehat{\mathbf{EDB}} \leftarrow \text{TWINSSE}_{\text{BASIC}}.\text{SETUP}(1^\lambda, \mathbf{DB})$$

$$\widehat{\mathcal{R}}_q = \text{TWINSSE}_{\text{BASIC}}.\text{SEARCH}(\text{sk}, \text{st}, q; \widehat{\mathbf{EDB}})$$

where $\mathcal{R}_q \subseteq \widehat{\mathcal{R}}_q$ and $|\widehat{\mathcal{R}}_q| \leq \frac{1}{\eta} \cdot |\mathcal{R}_q|$ ($0.85 < \eta \leq 1$) given that $\widehat{\mathcal{R}}_q$ is returned by $\text{TWINSSE}_{\text{BASIC}}$ (or $\mathbf{DB}(q_{\text{mkw}})$) and $\mathcal{R}_q = \mathbf{DB}(q)$.

Note that, the lower bound of η is obtained empirically from experiments over real databases. We present such experimental details in Section 5. The lower bound can be adjusted to accommodate larger $\widehat{\mathcal{R}}_q$ for different databases if required.

Storage Overhead. $\text{TWINSSE}_{\text{BASIC}}$ incurs $O(N^2)$ storage overhead to store the meta-keywords, where N is the number of keywords in the original plaintext database. This follows immediately from the fact that the number of meta-keywords is $O(N^2)$. This is undesirable in practice as it affects the scalability of the construction for large real databases. Currently, the schemes designed for disjunctive queries (such as IEX) require quadratic storage often leading to storage blow-up for large databases. Our final construction in Section 4 reduces $O(N^2)$ storage overhead to $O(N)$ - a necessary and significant reduction to use large real databases for deployment.

Search Overhead. The disjunctive search uses a meta-keyword as the least-frequent term for searching with the CSSE search routine. Since each mkw is an “union” of constituent ws, on average

⁶An exact solution returns only the documents belonging to the actual query result.

the frequency of the least-frequent mkw is smaller compared to a conjunctive query constituting the same ws. As a result, this basic method potentially can result to worst-case linear search overhead, which would be highly undesirable. However, we avoid such overheads by choosing an underlying CSSE scheme that ensures *sub-linear search complexity*.

4 TWINSSE: FINAL VERSION

In this section, we present our final scheme – TWINSSE, which improves upon $\text{TWINSSE}_{\text{BASIC}}$ with respect to storage requirements as well as search overheads. At the core of both these improvements lies an additional technique that we describe next – “frequency-based bucketization” of keywords. We note that similar techniques have been used in the SSE literature [17], albeit almost entirely for *frequency padding* and leakage-reduction. To the best of our knowledge, we are the first to show that bucketization can also be used to reduce storage and search overheads in SSE schemes.

4.1 Keyword Bucketization at Setup

We now describe our strategy for frequency-based keyword bucketization and intra-bucket meta-keyword generation at setup. We then use this updated meta-keyword generation strategy to formally describe the new setup algorithm – TWINSSE.SETUP .

Bucketization. Let $\Delta = \{w_1, \dots, w_N\}$ be a dictionary of keywords, and assume that these keywords are arranged in increasing order of frequency. Also, let $n' = O(1)$ be any arbitrarily chosen *constant*. We partition the keyword space into $n_B = N/n'$ “buckets” of size n' each ($n_B = \lceil \frac{N}{n'} \rceil$, if N is not a multiple of n'), where the k^{th} bucket is defined formally as the keyword subset

$$\Delta_k = \{w_{(k-1)n'+1}, \dots, w_{kn'}\}.$$

Note that since all keywords are arranged in increasing order of frequency, each bucket from Δ_1 through Δ_{n_B} progressively consists of keywords with increasing frequency ranges. We note that this is similar to the bucketization strategy employed in [17].

Intra-Bucket Meta-Keyword Generation. Having partitioned the keyword space into frequency-based buckets, we now proceed as follows:

- For each bucket Δ_k , we generate an *intra-bucket* meta-keyword set $\mathcal{S}_{\text{mkw},k}$ of size $O(|\Delta_k|^2) = O((n')^2)$. This is done exactly as in $\text{TWINSSE}_{\text{BASIC}}$, i.e., following the meta-keyword generation strategy in Lemma 3.1.
- We then define the overall set of meta-keywords as the collection of intra-bucket meta-keywords from all buckets, i.e.,

$$\mathcal{S}_{\text{mkw},\Delta} := \bigcup_{k \in [n_B]} \mathcal{S}_{\text{mkw},k}.$$

Observe that

$$|\mathcal{S}_{\text{mkw},\Delta}| = \sum_{k \in [n_B]} |\mathcal{S}_{\text{mkw},k}| = O(n_B(n')^2) = O(Nn').$$

However, $n' = O(1)$ is a constant, and hence, unlike in the basic solution described in Section 3, now $|\mathcal{S}_{\text{mkw},\Delta}| = O(N)$. In other words, we now have a linear-sized meta-keyword set, which forms the key stepping stone towards avoiding a quadratic storage overhead. We design our proposed TWINSSE to work for any choice

Algorithm 1 TWINSSE.SETUP

Input: $\text{DB}, 1^\lambda, n', n_B$
Output: $\text{sk}, \text{st}, \text{EDB}$

- 1: **function** TWINSSE.SETUP($1^\lambda, n', \text{DB}$)
- 2: $\widehat{\text{DB}} = \text{GENMETADB}(\text{DB}, n', n_B)$
- 3: $\text{sk}, \text{st}, \widehat{\text{EDB}} \leftarrow \text{CSSE.SETUP}(\lambda, \widehat{\text{DB}})$
- 4: **return** $\widehat{\text{EDB}}, \text{sk}, \text{st}$
- 5: Server receives $\widehat{\text{EDB}}$
- 6: Client keeps $(\text{sk}, \text{st}, n_B)$

of n' (ideally, n' should be a small constant to avoid high storage overheads); we use $n' = 10$ for our prototype implementation and experimentation over real-world databases in Section 5. We present brief discussion and empirical evaluations on the Enron dataset in Section 5 to select a suitable value for n' .

TWINSSE.SETUP: We now put these ideas together to formally describe TWINSSE.SETUP in Algorithm 1, which in turn again uses any generic conjunctive SSE scheme

$$\text{CSSE} = (\text{CSSE.SETUP}, \text{CSSE.SEARCH})$$

in a black-box way. The key changes from the basic scheme in Section 3 are highlighted in red for ease of exposition (in fact, TWINSSE_{BASIC}.SETUP can be viewed as a special case of TWINSSE.SETUP where all keywords are placed in the same bucket, i.e., $n_B = 1$ and $n' = N$).

Note that Algorithm 1 uses as a sub-routine Algorithm 2, which formally describes the meta-database generation based on the keyword bucketization and intra-bucket meta-keyword generation procedures described earlier. Overall, the working of Algorithm 1 can be divided into two steps: (a) generate the meta-database with the intra-bucket meta-keywords using Algorithm 2, and (b) generate the client state and the encrypted meta-database using CSSE.SETUP in a black-box way (note that this second step is the same as in TWINSSE_{BASIC}; the only alteration is in the generation of the meta-database, which now uses linearly many meta-keywords).

4.2 Updated Query Planning

We now describe the updated query planning strategy that takes into account the above mentioned meta-keyword generation process. We use this updated query planning strategy to build TWINSSE.SEARCH routine (the query planning for conjunctive queries remains same as in TWINSSE_{BASIC}).

At a high level, we partition a disjunctive query into “regions”, where each region consists of the keywords in the query that belong to the same bucket. Formally, given a query $q = (w_{\ell_1} \vee \dots \vee w_{\ell_n})$, let $Q = (w_{\ell_1}, \dots, w_{\ell_n})$ and, for each $k \in [n_B]$, let $Q_k = \Delta_k \cap Q$. In other words, Q_k consists of all the keywords in the disjunction q that belong to the k^{th} bucket. It is easy to see that we can re-write q as a disjunction over sub-queries as follows:

$$q = \bigvee_{k \in [n_B]} \left(\bigvee_{w \in Q_k} w \right) := \bigvee_{k \in [n_B]} q_k.$$

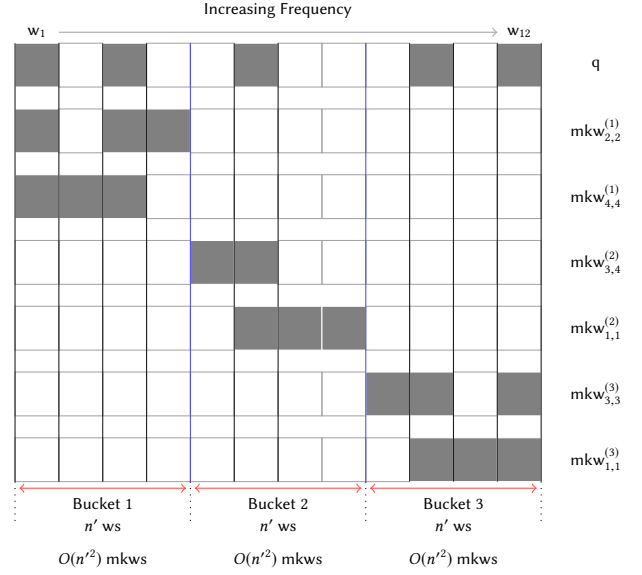


Figure 2: Expressing disjunctive query q in terms of mkws in the improved $O(N)$ construction. We use the same database parameters from Figure 1. The sub-mkws (smaller mkws in each bucket) are padded to realise a full mkw and stored in $\widehat{\text{DB}}$. Here, the bucket size n' is four. The k in $\text{mkw}_{i,j}^{(k)}$ represents the bucket index the $\text{mkw}_{i,j}$ belongs to. Note that, ANDing mkws within a bucket retains only the ws in the q (covered by that bucket). Note that expressing each meta-keyword as a bit-string allows efficient a transformation from the original disjunctive query into the corresponding conjunction of meta-keywords through simple bit-wise set/clear operations. This transformation incurs only negligible additional computational overhead in practice.

Note that for some k , Q_k could be an empty set; in this case, the sub-query q_k is also empty. Based on the above representation, the query planning strategy in TWINSSE works as follows:

- **Step-1:** Partition q into sub-queries $\{q_k\}_{k \in [n_B]}$ as described above.
- **Step-2:** For each sub-query q_k , construct a conjunctive (sub-) meta-query $q_{\text{mkw},k}$ as described in Section 3, using the intra-bucket meta-keywords corresponding to the k^{th} bucket, i.e., the intra-bucket meta-keywords in $S_{\text{mkw},k}$.
- **Step-3:** Finally, we define the overall meta-query q_{mkw} as $q_{\text{mkw}} = \bigvee_{k \in [n_B]} q_{\text{mkw},k}$, and re-construct $\text{DB}(q_{\text{mkw}})$ as

$$\text{DB}(q_{\text{mkw}}) = \bigcup_{k \in [n_B]} \text{DB}(q_{\text{mkw},k}),$$

where the recovery of each $\text{DB}(q_{\text{mkw},k})$ happens via an independent (and parallel) execution of the same search protocol as in TWINSSE_{BASIC}.

The above query planning strategy is summarized pictorially in Figure 2 (note that the superscript k in $\text{mkw}_{i,j}^{(k)}$ represents the bucket index for the meta-keyword $\text{mkw}_{i,j}$). In comparison with

Algorithm 2 GENMETADB

Input: DB, n', n_B
Output: $\widehat{\text{DB}}$

```

1: function GENMETADB( $\text{DB}, n', n_B$ )
2:   Extract  $\Delta$  from  $\text{DB}$  and sort ws in  $\Delta$  in increasing order of
   frequency
3:   Partition  $\Delta$  as  $\{\Delta_1, \dots, \Delta_{n_B}\}$  such that  $\Delta_i =$ 
    $\{w_{(i-1)n'+1}, \dots, w_{in'}\}$   $\triangleright$  The last bin may not contain  $n'$ 
   ws. Keep only as many are left.
4:   Initialise a bucket index  $k \leftarrow 1$ 
5:   for  $\Delta_k \in \{\Delta_1, \dots, \Delta_{n_B}\}$  do
6:     Set  $\ell \leftarrow |\Delta_k|$ 
7:     Parse  $\Delta_k$  as  $\{w_1^k, \dots, w_\ell^k\}$ 
8:     for  $i \leftarrow 1$  to  $\ell$  do
9:       for  $j \leftarrow 0$  to  $\ell - i$  do
10:         $\text{mkw}_{i,i+j}^{(k)} \leftarrow \Delta_k \setminus \{w_i^k, \dots, w_{i+j}^k\}$ 
11:        for  $w \in \text{mkw}_{i,i+j}^{(k)}$  do
12:           $\widehat{\text{DB}}(\text{mkw}_{i,i+j}^{(k)}) \leftarrow \widehat{\text{DB}}(\text{mkw}_{i,i+j}^{(k)}) \cup \text{DB}(w)$ 
13:         $k \leftarrow k + 1$ 
14:   return  $\widehat{\text{DB}}$ 

```

Algorithm 3 TWINSSE.SEARCH (for disjunctive queries)

Input: $q, \text{sk}, \text{st}, \widehat{\text{EDB}}, n', n_B$
Output: Result set $\text{DB}(q)$

```

1: function TWINSSE.SEARCH( $q, \text{sk}, \text{st}, \widehat{\text{EDB}}$ )
2:   Client
3:   Generate  $q_{\text{mkw}} = (\bigvee_{k \in [n_B]} q_{\text{mkw},k}) =$ 
   GENMQUERY( $q, n', n_B, \Delta$ ).
4:   For each non-empty  $q_k$  (in uniformly random order), the
   client and server engage in the search protocol as below.
5:   Client+Server
6:   for each non-empty  $q_{\text{mkw},k}$  (in random order) do
7:      $\text{DB}(q_{\text{mkw},k}) \leftarrow$ 
   CSSE.SEARCH( $\text{sk}, \text{st}, q_{\text{mkw},k}; \widehat{\text{EDB}}$ ).
8:     At the end of the protocol, client receives  $\text{DB}(q_{\text{mkw},k})$ .
9:   Client
10:  Initialize  $\text{DB}(q_{\text{mkw}}) \leftarrow \text{EMPTY-SET}$ .
11:  for each  $\text{DB}(q_{\text{mkw},k})$  from search protocol do
12:     $\text{DB}(q_{\text{mkw}}) \leftarrow \text{DB}(q_{\text{mkw}}) \cup \text{DB}(q_{\text{mkw},k})$ .
13:  The client locally filters  $\text{DB}(q) \subseteq \text{DB}(q_{\text{mkw}})$ .

```

the example figure (Figure 1) in Section 3, we note that the meta-keywords are now chosen from a smaller set of size $\approx (4 \times 12) = 48$, as compared to a set of size $\approx 12^2 = 144$ in Figure 1.

TWINSSE.SEARCH: We now put these ideas together to formally describe TWINSSE.SEARCH in Algorithm 3, which in turn uses Algorithm 4 as a sub-routine (we only summarize the processing of disjunctive queries since conjunctive queries are processed as in TWINSSE_{BASIC}). The key changes from the basic scheme in Section 3 are highlighted in red for ease of exposition (again, TWINSSE_{BASIC}.SEARCH can be viewed as a special case of TWINSSE.SEARCH where all keywords are placed in the same bucket, i.e., $n_B = 1$ and $n' = N$).

Algorithm 4 GENMQUERY

Input: $q = \{w_{i_1} \vee w_{i_2} \vee \dots \vee w_{i_l}\}, n', n_B, \Delta$
Output: Meta-Query q_{mkw}

```

1: function GENMQUERY( $q, n', n_B, \Delta$ )
2:   Parse  $\Delta$  as  $\{\Delta_1, \dots, \Delta_{n_B}\}$ 
3:   Sort ws in  $q$  in increasing order of frequency
4:   Partition query  $q$  into set of sub-queries  $P_q$  as  $q_1 || \dots || q_{n_B}$ ,
   such that  $q_k$  contains ws only from  $\Delta_k$  for  $k = 1, \dots, n_B$ 
5:   for  $q_k \in P_q$  do
6:     Parse  $q_k$  as  $\{w_{i_1}^k, \dots, w_{i_{l'}}^k\}$ 
7:     for  $j \leftarrow 1$  to  $l'$  do
8:        $\text{mkw}_{i_{j-1}+1, i_j-1}^{(k)} \leftarrow \Delta_k \setminus \{w_{i_{j-1}+1}^k, \dots, w_{i_j-1}^k\}$ 
        $\triangleright$  Recall that  $\text{mkw}_{i_k, j_k} \leftarrow \phi$  if  $i_k > j_k, i_0 = 0$  and
        $w_1^k$  is the first keyword in  $\Delta_k$ 
9:        $\text{mkw}_{i_{l'}+1, n'}^{(k)} \leftarrow \Delta_k \setminus \{w_{i_{l'}+1}^k, \dots, w_{n'}^k\}$ 
10:       $q_{\text{mkw},k} \leftarrow \text{mkw}_{1, i_1-1}^{(k)} \wedge \dots \wedge \text{mkw}_{i_{l'}+1, n'}^{(k)}$ 
11:       $q_{\text{mkw}} \leftarrow q_{\text{mkw}} \vee q_{\text{mkw},k}$ 
12:   return  $q_{\text{mkw}}$ 

```

Algorithm 4 formally captures the updated disjunctive query planning strategy based on query partitioning and intra-bucket meta-keywords, as described earlier. Note that in Algorithm 3, each conjunctive sub-meta-query q_k is executed in parallel using the search algorithm CSSE.SEARCH of the underlying conjunctive SSE scheme CSSE, and the final result-set corresponding to the overall meta-query is constructed locally at the client by taking the union over the result-sets corresponding to each conjunctive sub-meta-query.

Correctness. We state the following theorem for the correctness of TWINSSE.

THEOREM 4.1 (CORRECTNESS OF TWINSSE). *Assuming that CSSE satisfies correctness of search for conjunctive queries and Lemma 3.1 holds, TWINSSE satisfies correctness for both conjunctive and disjunctive queries.*

The proof essentially follows from the same arguments as the proof of correctness for TWINSSE_{BASIC} in Section 3 and is presented in Appendix B.

4.3 Instantiation from the OXT Protocol and Complexity Analysis

In its most general form, our proposed TWINSSE scheme can be concretely instantiated using any conjunctive SSE scheme. In this section, we analyze a concrete instance of TWINSSE based on the OXT protocol [7], which we call TWINSSE_{OXT}. We analyze TWINSSE_{OXT} asymptotically in terms of storage requirements and search overheads. Our analysis does not require understanding the internal details of OXT beyond what is already stated in this section; the reader may refer to [7] for more details. Finally, we refer the reader to Section 5 for experimental validation of the analysis presented here over the Enron email corpus.

Storage Requirements (Server). The (worst-case) server-side storage requirement for TWINSSE_{OXT} is $O(n'|\text{DB}|)$, where $|\text{DB}|$ is the number of distinct identifier-keyword pairs in DB , and $n' = O(1)$

denotes the (constant) size of each keyword bucket used. This linearization process through keyword bucketization process incurs an $O(n')$ -fold increase in storage overhead over OXT (where $n' = O(1)$ is a constant). We view this as a necessary trade-off for the additional ability to support disjunctive queries efficiently yet securely. In comparison, the IEX family of schemes incur (worst-case) quadratic storage overheads, more precisely, $O(|\Delta||\mathbf{DB}|)$, where $|\Delta|$ denotes the number of keywords in \mathbf{DB} .

However, n' (or the number of buckets) needs to be chosen carefully to bound the storage overhead to linear (which also keeps the leakage from multiple buckets at minimum). A high value of n' would incur a higher storage overhead with lesser leakage from small number of buckets (as outlined in Section 4.1). Whereas a small value of n' would result in a higher number of buckets leading to lesser storage but increased leakage from more number of buckets. We selected n' in the range 10-15 based on empirical evaluations over real data sets that allows to retain a linear storage overhead. These experimental results are provided in Section 5.

Storage Requirements (Client). The client-side storage requirement for $\text{TWINSSE}_{\text{OXT}}$ is $O(n'|\Delta| \log |\mathbf{DB}|)$. This again represents an $O(n')$ -fold increase in storage overhead over the original OXT scheme (where $n' = O(1)$ is a constant), which has a client-side storage requirement of $O(|\Delta| \log |\mathbf{DB}|)$.

We also note here that $\text{TWINSSE}_{\text{OXT}}$ requires $O(1)$ storage for the secret key(s) at the client-end (this is a purely client-side overhead, not associated with the server-side storage). In contrast, the IEX family of schemes require $O(|\Delta|)$ secure storage for secret keys (one key per keyword due to individual multi-map structure required for each keyword index), which is likely to be costly for extremely large databases. We emphasise that this requirement is only for *secure* storage to store secret keys on the client-side. The client-side storage overhead mentioned in the previous paragraph accounts for storing only auxiliary information required during query processing, and this *does not* require secure storage.

Search Complexity. We now present an asymptotic analysis of the search complexity (more concretely, the computational and communication requirements during search query processing) of $\text{TWINSSE}_{\text{OXT}}$. We divide our analysis into two parts – conjunctive queries and disjunctive queries:

Conjunctive queries. Let $q = (w_1 \wedge \dots \wedge w_n)$ be a conjunctive query, where w_1 is the *least frequent* keyword. When processing q using $\text{TWINSSE}_{\text{OXT}}$, the computational costs (at both the client and the server) as well as the communication requirements between the client and the server scale linearly as $O(n|\mathbf{DB}(w_1)|)$. This is *exactly the same* as in OXT, and is hence worst case sub-linear.

Disjunctive queries. Let $q = (w_1 \vee \dots \vee w_n)$ be a disjunctive query. Also, let $q_{\text{mkw}} = \vee_{k \in [n_B]} q_{\text{mkw},k}$ be the corresponding meta-query, and assume without loss of generality that $\text{mkw}_{i_k, j_k}^{(k)}$ is the least frequent meta-keyword within $q_{\text{mkw},k}$ for each $k \in [n_B]$ (such that $q_{\text{mkw},k}$ is non-empty). When processing q using $\text{TWINSSE}_{\text{OXT}}$, the computational costs (at both the client and the server) as well as the communication requirements between the client and the server scale linearly as $O(\gamma)$, where

$$\gamma = \sum_{k \in [n_B]} |q_k| |\mathbf{DB}(\text{mkw}_{i_k, j_k}^{(k)})|,$$

where $|q_k|$ denotes the number of meta-keywords in the conjunctive sub-meta-query q_k ($|q_k| = 0$ when q_k is empty). Note that this is essentially a generalization of the analysis of search query overheads for $\text{TWINSSE}_{\text{BASIC}}$ in Section 3, where all keywords belong to the same bucket (i.e., $n_B = 1$). We provide a comparative summary of storage and search overhead for TWINSSE and IEX in Table 1 for quick reference.

Table 1: Comparative summary of storage overhead and search complexity of TWINSSE and IEX.

Scheme	Storage Overhead	Search Time ⁷	
		Conjunctive	Disjunctive
$\text{TWINSSE}_{\text{OXT}}$	$O(n' \mathbf{DB})$	$O(n \mathbf{DB}(w_1))$	$\sum_{k \in [n_B]} q_k \mathbf{DB}(\text{mkw}_{i_k, j_k}^{(k)}) $
IEX-2LEV	$O(\Delta \mathbf{DB})$	$O(n^2(\mathbf{DB}(w_M^u) + t \mathbf{DB}(\delta_u)))$	$O(n^2 \mathbf{DB}(w_M))$

Spurious document identifiers. It turns out that keyword bucketization also significantly reduces the search overheads (both computational and communication) due to spurious document identifiers in $\mathbf{DB}(q_{\text{mkw}})$. In particular, recall our observation with respect to $\text{TWINSSE}_{\text{BASIC}}$ from Section 3: the fraction of spurious identifiers retrieved is directly proportional to the *average number of common documents* over every keyword-pair in the database. However, in our improved solution, the database is partitioned into buckets, and all keywords within the same bucket have essentially similar frequency ranges. This means, in particular, that an overwhelmingly large fraction of buckets either contain all low-frequency keywords (in which case, the spurious document-set is essentially null, since such keywords almost never co-occur across documents [7]), or very high-frequency keywords (in which case, such keywords occur in almost all documents, and the proportion of spurious documents is low by default).

We generalise the aforementioned observations into the following (informal) claim about the search complexity incurred by $\text{TWINSSE}_{\text{OXT}}$, which is essentially an extension of our claim for $\text{TWINSSE}_{\text{BASIC}}$.

CLAIM 4.1 (INFORMAL). $\text{TWINSSE}_{\text{OXT}}$ incurs (average-case) sub-linear search complexity (in terms of both computational costs and communication overheads) for both conjunctive and disjunctive queries.

We validate this claim with experimental results over the Enron email corpus in Section 5. We also extend the analysis and experimental evaluations to more general/complicated Boolean queries (CNF or DNF formulae) in Appendix G. Our experiments show that searches in $\text{TWINSSE}_{\text{OXT}}$ incur at most 15% overhead due to spurious identifiers in the result set. In order to filter out the set of spurious documents from the final result set, we can resort to the same strategies as used by state-of-the-art volume hiding SSE constructions (e.g. SSE schemes obtained naturally from the encrypted multi-map constructions proposed in [20, 28], where the client obtains a mixture of real and “fake” identifiers at the end of the query phase to hide the true query response volume from the server).

⁷ w_1 is the least frequent keyword in the conjunctive query, w_M^u is the most frequent keyword in the CNF clause δ_u with most number of w s in CNF expression of the conjunctive query. w_M is the most frequent keyword in the disjunctive query.

Note that, TWINSSE uses separate mkw^* for each bucket where mkw_i^* denotes the mkw^* for the i -th bucket. Following the definition of mkw^* , mkw_i^* represents the disjunction of all w s in the i -th bucket. Using a separate mkw^* per bucket allows TWINSSE to support non-SNF queries more efficiently than OXT. OXT uses a single mkw^* for any non-SNF query into an SNF query. In this process, OXT incurs a worst-case linear search overhead. In contrast, TWINSSE uses mkw_i^* if and only if a query involves *all* keywords from the i -th bucket (in this case, it is the only optimal choice). These specific queries can be considered as corner cases which rarely occur in realistic searches. Consequently, TWINSSE incurs spurious ids, which is typically only 15 – 20% on average for realistic non-SNF queries. Additionally, this process in TWINSSE_{OXT} also allows for parallel execution of independent sub-queries over different buckets. Adopting a similar approach with OXT for parallel execution would incur more leakage from each queried bucket due to the exact result set which reveals the volume pattern. Since TWINSSE_{OXT} produces noisy result set due to the spurious ids, the volume pattern leakage is less compared to OXT.

4.4 Security of TWINSSE

We present an informal discussion on the security of our construction here. Detailed formal security discussion and leakage analysis (including experimental evaluations) are available in Appendix C, D, E and F. Security of TWINSSE is modelled in the semi-honest adversarial setting where the server is assumed to be a honest-but-curious entity (that means, the server follows the algorithmic specifications exactly, but can record information for later analysis).

Informally, TWINSSE inherits security properties and leakage profile from the underlying CSSE construction. We assume that the underlying CSSE construction is an adaptively secure sublinear conjunctive SSE algorithm which is secure against a semi-honest adversary \mathcal{A} and the leakage of CSSE is characterised by the leakage function \mathcal{L}_{CSSE} . The leakage function \mathcal{L}_{CSSE} is an ensemble of the leakage functions for **SETUP** and **SEARCH** individually, expressed in the following way.

$$\mathcal{L}_{CSSE} = \{\mathcal{L}_{CSSE}^{SETUP}, \mathcal{L}_{CSSE}^{SEARCH}\}$$

Given the above CSSE leakage functions, security of TWINSSE can be analysed using TWINSSE leakage function $\mathcal{L}_{TWINSSE}$ in the same adaptive semi-honest adversarial model. Similar to \mathcal{L}_{CSSE} , $\mathcal{L}_{TWINSSE}$ is composed of two separate leakage functions for **SETUP** and **SEARCH**, as expressed below, that capture the leakage from TWINSSE execution in the meta-keyword setting.

$$\mathcal{L}_{TWINSSE} = \{\mathcal{L}_{TWINSSE}^{SETUP}, \mathcal{L}_{TWINSSE}^{SEARCH}\}$$

Concretely, $\mathcal{L}_{TWINSSE}$ is identical to the \mathcal{L}_{CSSE} with n_B (the number of buckets) as an additional benign component. In other words, we show that $\mathcal{L}_{TWINSSE}$ is equal to $\tilde{\mathcal{L}}_{CSSE}$ where $\tilde{\mathcal{L}}_{CSSE}$ is \mathcal{L}_{CSSE} in the context of meta-keywords and n_B . At a high level, $\mathcal{L}_{TWINSSE}^{SETUP}$ incorporates \widehat{DB} instead of DB generated by the GENMETADB during setup. Similarly, the search leakage encapsulates leakages from *both* conjunctive and disjunctive queries. We quantify these separately through two individual leakage function instances - one for conjunctive queries, and one for disjunctive queries from meta-keywords, where meta-keywords are generated using GENMQUERY routine. We show that the leakage for the conjunctive case is exactly

the same as of the CSSE construction, and for disjunctive queries it incorporates the similar leakage profile, but from meta-keywords. We provide detailed formal analysis of $\mathcal{L}_{TWINSSE}$ in Appendix D.

Security of TWINSSE_{OXT}. The security analysis of TWINSSE_{OXT} follows from the security notions of generic TWINSSE, as informally discussed above (formally in Appendix D and Appendix E). Due to lack of space, we move this discussion (including proofs) to the Appendix D. We also present a leakage-based cryptanalysis of the TWINSSE_{OXT} scheme via experiments over the Enron email corpus in Appendix F.

Comparison with IEX. At a high level, TWINSSE_{OXT} avoids two kinds of leakages that the IEX family of schemes incurs for any query. To begin with, IEX leaks to the server the exact size of the result set pertaining to a query (also referred to as the size pattern leakage). As already mentioned, due to the presence of spurious document identifiers in the result set, TWINSSE_{OXT} inherently hides the size pattern from the server. More crucially, IEX incurs significant sub-query leakage. For example, given a disjunctive query of the form $q = (w_1 \vee w_2)$, where w_1 is the more frequent keyword, it leaks to the server: (a) the frequency of the more frequent keyword, i.e., $|\mathbf{DB}(w_1)|$, and (b) the number of documents that contain w_2 but not w_1 , i.e., $|\mathbf{DB}(w_2) \setminus \mathbf{DB}(w_1)|$. Whereas, TWINSSE_{OXT} only leaks the frequency of the least frequent meta-keyword (in this example, the meta-keyword corresponding to w_1), and no information about the other meta-keywords in the conjunction (in this example, no information about w_2). In other words, TWINSSE_{OXT} incurs less leakage as compared to the IEX family of schemes during search queries.

5 EXPERIMENTAL RESULTS

We describe a prototype implementation of TWINSSE_{OXT} and evaluate its performance over real-world databases. We present experimental results comparing the storage requirements and search performance of TWINSSE_{OXT} with that of IEX-2LEV [19].

Data Set and Platform. We used the Enron email data set^{8,9} for our experiments. The Enron email data set contained 517,401 documents (emails) and 20 million keyword-document pairs, with a total size 1.9 GB. The complete TWINSSE_{OXT} implementation was done using C++ (with GCC9 compiler) with native multi-threading support, and we used Redis as the database backend. We ran the experiments on a *single* node with Intel Xeon E5-2690 v4 2.6 GHz CPU with 128 GB RAM and 512 GB SSD storage.

Implementation Details. We created the meta-keyword database (or the transformed database) **MDB** from the parsed Enron database **DB**. The plain Enron database **DB** contains w -s and id -s in inverted index form. The transformed database **MDB** also contains the mkw s and the associated id -s in inverted index layout. Since there are a large number of w -s in **DB**, length of each binary string mkw is large. Hence, we hash those strings prior to writing to **MDB**. This **MDB** is further encrypted using the underlying OXT setup to generate the encrypted meta-keyword database **EDB**.

We report the actual size of **EDB** in Figure 3 which is offloaded to the server. The query translation process first generates these

⁸<https://www.cs.cmu.edu/~enron/>

⁹<https://www.kaggle.com/wcukierski/enron-email-dataset>

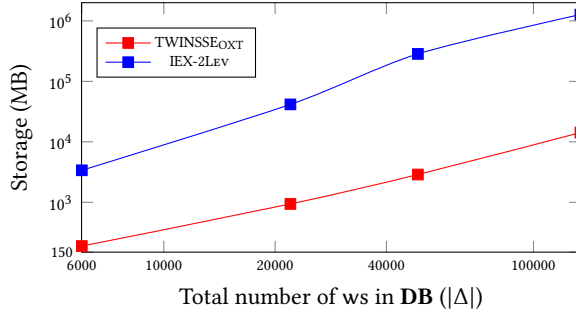


Figure 3: Server storage overhead with database size ($|\text{DB}|$) for the Enron database.

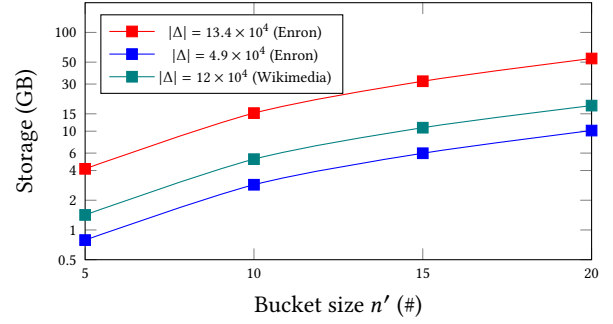


Figure 5: Variation of storage overhead for different choices of the bucket size (n').

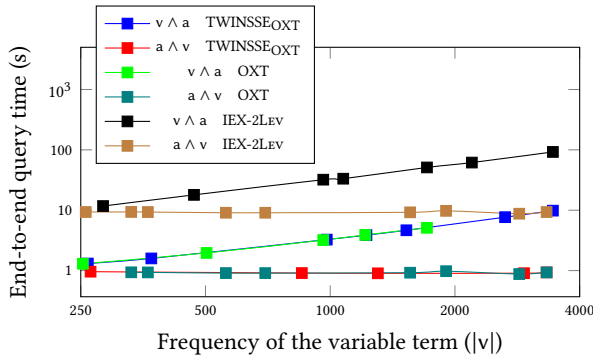


Figure 4: Comparison of end-to-end search latency vs frequency of the variable term ($|v|$). Observe that, for conjunctive queries TWINSSE_{OXT} closely follows OXT in practice, and the latency is significantly less than IEX-2LEV. We note here that a fundamental difference between IEX-2LEV and TWINSSE_{OXT} is that the search frequency of IEX-2LEV scales (by design) with the frequency of the most frequent conjunct, while that of TWINSSE_{OXT} scales with the frequency of the least frequent conjunct; this is the main reason why TWINSSE_{OXT} outperforms IEX-2LEV by a significant margin for conjunctive queries.

mkws in binary string format and we hash those prior to search over the encrypted meta-keyword database $\widehat{\text{EDB}}$.

Evaluation of Storage Overhead. One of the fundamental aspects of our implementation is that TWINSSE_{OXT} improves upon the quadratic storage overhead of IEX-2LEV and scales linearly with the size of plaintext database. IEX-2LEV exploits the low size of mutual intersections for all pairs of keywords in DB and its storage overhead scales with the size of the intersection. In a sparse data set, the size of these intersections for most of the pairs of keywords is very low. However, if the database is not sparse, this results in large intersections for pairs of w -s and the overhead becomes truly quadratic for IEX-2LEV.

Figure 3 compares the storage overhead of TWINSSE_{OXT} and IEX-2LEV on Enron database (sparse database). It is evident that the storage size scales linearly with the number of keywords in DB

for TWINSSE_{OXT} whereas IEX-2LEV becomes quadratic leading to storage blow-up. The storage overhead of IEX-2LEV is 60× more than TWINSSE_{OXT}. Despite the additional storage required for the meta-keywords, TWINSSE_{OXT} has better storage overhead in worst-case distribution of DB as compared to IEX-2LEV.

Effect of linearization. As discussed in Section 4.3, the choice of n' greatly influences the storage overhead. Since the distribution of w s and ids varies across different databases (for example, a medical database’s distribution differs from a tax record database), it is quite challenging (and inefficient) to obtain an analytic expression for n' that works for multiple databases. We rely on an empirically chosen value of n' that suitably works for different databases without blowing up the storage. We present experimental results in Figure 5 to illustrate the effect of varying n' . We fix n' at 10 for our final experiments from this evaluation.

Evaluation of End-to-End Search Latency. Figure 4 and 6 compare the end-to-end search latency of TWINSSE_{OXT} with that of IEX-2LEV for conjunctive and disjunctive queries, respectively.

Conjunctive queries. Search performance of TWINSSE_{OXT} is inherited from OXT and is therefore identical to OXT as shown in Figure 4. To validate this, we consider a two-keyword query of the form $q = a \wedge v$, where a and v are two keywords from DB . Without loss of generality, we consider the first term of q (or a here) to be the least frequent keyword. We vary the frequency of v (referred to as the variable term) with different queries where as the frequency of a is kept constant (constant term). The plot shows constant time overhead for conjunctive queries of this form with TWINSSE_{OXT}, which is identical with OXT. In the same figure, IEX-2LEV conjunctive search time is also plotted which depicts that TWINSSE_{OXT} is around 10× faster on average.

Disjunctive queries. The plots in Figure 6 compare the end-to-end query time with final result size for disjunctive queries of different hamming weights. Observe that, the query time increases with increasing number of id -s in $\widehat{\mathcal{R}}_q$ (the obtained result set, inclusive of the spurious id -s for a disjunctive query) due to the increased frequency of least-frequent mkws in these queries.

As discussed in the main text, frequency of the least-frequent mkw is independent of the frequency of the least-frequent w in a query. Hence, we consider plotting with overall result size that represents the computation overhead. In disjunctive queries, union of

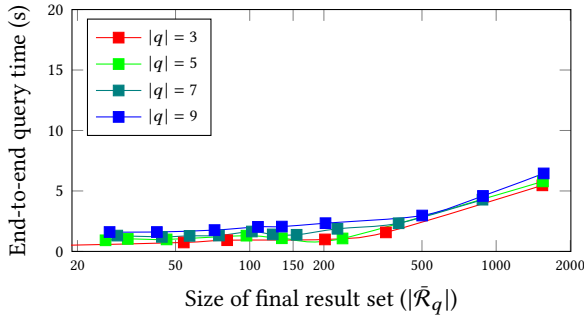


Figure 6: End-to-end search latency vs final result size for disjunctive queries of different number of ws in q ($|q|$).

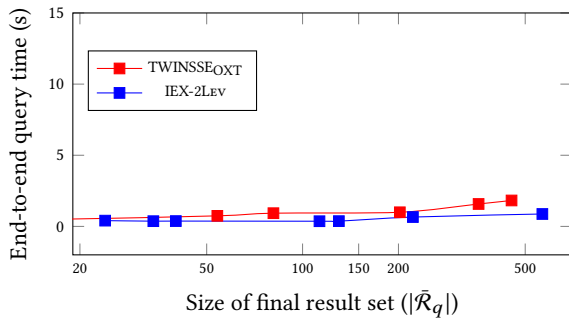


Figure 7: Comparison of end-to-end search latency vs final result size for disjunctive queries of the form $q = w_1 \vee w_2$ over Enron data set.

the id-s grows with more number of keywords present in the query. Therefore, plotting with the result size provides an accurate measure of computation cost for disjunctive queries. Nonetheless, the OXT sublinear search complexity is maintained, which we verified in our experiments.

The average query time increases with the number of ws in actual disjunctive query q . This increased time attributes to more number of mkws for each query, and the underlying OXT that scales linearly with number of keywords (in this case mkws) in the conjunctive query. The end-to-end disjunctive search latency for TWINSSE_{OXT} is few hundred milliseconds over the Enron database for queries with moderate result size.

We provide an end-to-end query performance comparison of TWINSSE_{OXT} with IEX-2LEV in Figure 7. For queries with smaller result sizes, TWINSSE_{OXT} achieves almost identical end-to-end query latency as IEX-2LEV. For queries with larger result sizes, IEX-2LEV performs slightly better. This is primarily because of the usage of relatively costly elliptic-curve cryptography-based operations in TWINSSE_{OXT} (a consequence of using OXT as a black-box, which uses such operations); IEX-2LEV, on the other hand, uses purely symmetric-key crypto-primitives. We view this as an efficiency trade-off; note that TWINSSE_{OXT} outperforms IEX-2LEV significantly both in terms of storage requirements (as demonstrated in Figure 3) and end-to-end latency for conjunctive queries (as

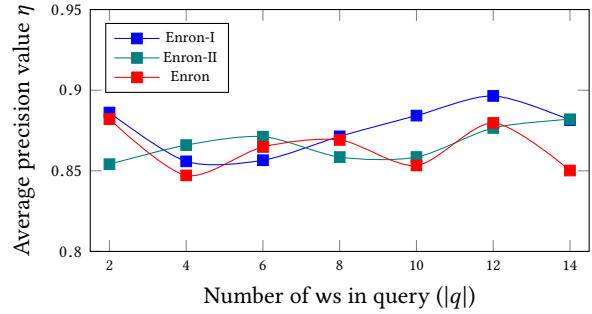


Figure 8: Average precision value vs number of keywords in query for databases of different size.

demonstrated in Figure 4). Hence, from the point of view of practical performance across a wide class of Boolean queries generally encountered in practice and scalability to extremely large databases, TWINSSE_{OXT} outperforms IEX-2LEV.

Experiments on the Wiki database. We also present experimental results for performance and storage overhead evaluation over the Wikimedia dataset¹⁰ in Appendix H. We observed similar results with the Wikimedia dataset as of the Enron dataset presented above.

Evaluation of Result Precision. In context of information retrieval *precision* (denoted by η in Section 3.3) is the fraction of relevant documents among the retrieved documents. We compare the average precision values of $\bar{\mathcal{R}}_q$ for disjunctive queries (q) with different number of keywords in Figure 8. Observe that the average precision values for most of the cases is above 85%, which implies that at least 85% documents returned by $\bar{\mathcal{R}}_q$ are relevant to the disjunctive query q (or belongs to the actual result set \mathcal{R}_q without spurious id-s). The plot also illustrates that scaling the database does not affect the average precision of the retrieved documents. Hence, the query result of TWINSSE_{OXT} does not degrade even for huge databases which is crucial for practical applications.

Note that, IEX is an exact solution that has 100% result precision - it returns the exact result set without spurious ids¹¹. However, IEX incurs extremely high storage overhead that makes it impossible to deploy with large real datasets. In contrast, TWINSSE incurs less than 100% result precision (85%-90%, as shown in Figure 8), but TWINSSE outweighs the loss in precision with storage savings ($10\times-50\times$ less than IEX, as shown in Figure 3).

Complex Boolean Queries. We defer the elaborate discussion on processing complex Boolean queries (as CNF and DNF formulae) using TWINSSE_{OXT} to Appendix G due to lack of space. We present here the corresponding experiments to evaluate the performance of our prototype implementation of TWINSSE_{OXT} when processing such queries over the Enron dataset (on the same computing platform as discussed in Section 5).

DNF queries. We considered multiple queries with two clauses and three clauses with each clause having two keywords. The end-to-end query time is plotted in Figure 15, where the blue curve

¹⁰<https://dumps.wikimedia.org/enwiki/latest/>

¹¹State-of-art SSE schemes like OXT[7], HXT[23], ODXT[29] or IEX[19] are exact solutions. Hence, this precision parameter is defined exclusively for TWINSSE only which produces result set with spurious ids

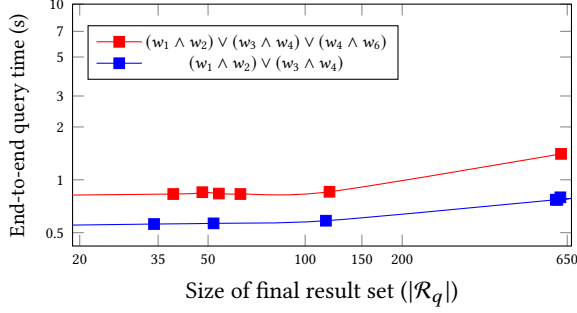


Figure 9: TWINSSE_{OXT} performance with result set size on Enron data set for DNF queries.

represents the query time for two clause queries and the red curve represents the query time for three clause queries. Observe that the query time for both two and three clause queries increase with more number of id-s in the final result set. This increment can be attributed to large result size of the individual conjunctive clauses. Also note that the query time increases for three-clause queries due to more conjunctive clauses and follows the same trend of increased query time with the final result size.

CNF queries. For experimenting with CNF queries, we considered two-clause queries with two and three keywords per clause. Since the Enron data set is relatively sparse in nature, it often results in small or empty intersection with higher number of clauses in query. We plotted end-to-end query time in Figure 16 for both cases – two keyword clauses and three keyword clauses with the size of the final result set. The blue curve represents the end-to-end query time for the queries with two keywords per clause. Similarly, the red curve represents the end-to-end query time for queries with three-keyword clauses. Observe that, in CNF queries also, the end-to-end query time increases with the final result size, due to the increased size of the initial result set retrieved. For the three-keyword clauses, the query time is higher than the two-keyword clauses due to the larger size of the initial result set retrieved for the disjunctive clauses.

Normally, CNF and DNF expressions in generic form involve the negation of terms. TWINSSE avoids negations for plain disjunctive and DNF queries. This is purposefully done to avoid quadratic storage due to negated keywords, as discussed in Section 1.2. However, TWINSSE_{OXT} supports negations in CNF by utilising the XSet data structure as a part of the underlying OXT construction. In particular, TWINSSE_{OXT} can process keyword negation for any clause in a CNF query, except for the least-frequent one (which is used to retrieve encrypted ids), by evaluating the formula directly over the XSet structure. Since the retrieval is done using the least frequent clause, TWINSSE_{OXT} avoids the inherent linear search overhead incurred by this kind of queries in OXT.

The authors of OXT [7] described a way to process CNF/DNF queries using a special keyword (similar to mkw^* in TWINSSE) as the *s-term* (special term or the keyword with minimum frequency in the query). This process essentially retrieves all entries from the database due to this special keyword and filters according to the query. Naturally, this approach has linear complexity that leads

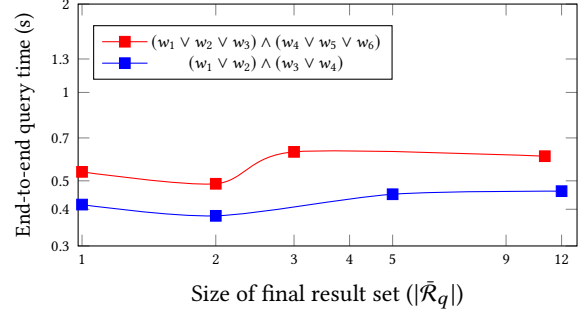


Figure 10: TWINSSE_{OXT} performance with result set size on Enron data set for CNF queries.

to poor performance and high leakage (a very high number of ids are retrieved, which do not belong to the final result). In contrast, TWINSSE_{OXT} still offers a sublinear search for complex queries. The all-1 mkw^* is defined per-bucket in TWINSSE_{OXT} scheme, which avoids such costly retrieval as only specific buckets are accessed and, consequently, leaks much less compared to OXT.

Storage Overhead for Dense Databases. Finally, we compare the server-side storage overheads of TWINSSE_{OXT} and IEX-2LEV for a special class of “dense databases” where for any pair of keywords, the number of documents containing both keywords is large (see Appendix I for an example of dense database and elaborate discussion). We use a synthetic dense database that follows Zipf’s law [24]. Our experimental results in Appendix I (Figure 20) show that IEX-2LEV incurs 70× higher storage overhead than TWINSSE_{OXT} for this dense database, clearly indicating that TWINSSE_{OXT} offers significantly greater scalability for such databases in practice.

6 SUPPORTING DYNAMIC DATABASES

In this paper, we described TWINSSE_{OXT} for static databases. This leaves open the question of extending TWINSSE_{OXT} to dynamic databases, and supporting updates efficiently yet securely over these. We note here that for dynamic databases where the set of keywords across all documents remains fixed (or, more generally, undergoes updates infrequently), the set of meta-keywords also does not change (frequently) over time. In this setting, it is possible to achieve an extension of TWINSSE_{OXT} to the setting of dynamic databases by simply substituting the underlying OXT scheme with a dynamic conjunctive SSE scheme with desirable efficiency and security guarantees (e.g. ODXT from [29]). However, such an extension becomes challenging for dynamic databases where the set of keywords (and hence, the set of meta-keywords) also gets updated frequently. We leave it as an interesting future question to extend TWINSSE_{OXT} to dynamic databases.

ACKNOWLEDGMENTS

The authors would like to thank the grant “Design and Implementation of Efficient and Secure Searchable Encryption” sponsored by MHRD-STARS (Scheme for Transformational and Advanced

Research in Sciences), India and Centre on Hardware Security Entrepreneurship Research & Development, Meity India for partially supporting the work.

REFERENCES

- [1] Ghous Amjad, Sarvar Patel, Giuseppe Persiano, Kevin Yeo, and Moti Yung. 2021. Dynamic Volume-Hiding Encrypted Multi-Maps with Applications to Searchable Encryption. *IACR Cryptol. ePrint Arch.* (2021), 765.
- [2] Laura Blackstone, Seny Kamara, and Tarik Moataz. 2020. Revisiting Leakage Abuse Attacks. In *NDSS 2020*.
- [3] Raphael Bost. 2016. $\Sigma\sigma\phi\sigma$: Forward Secure Searchable Encryption. In *ACM CCS 2016*. 1143–1154.
- [4] Raphaël Bost, Brice Minaud, and Olga Ohrimenko. 2017. Forward and Backward Private Searchable Encryption from Constrained Cryptographic Primitives. In *ACM CCS 2017*. 1465–1482.
- [5] David Cash, Paul Grubbs, Jason Perry, and Thomas Ristenpart. 2015. Leakage-Abuse Attacks Against Searchable Encryption. In *ACM CCS 2015*. 668–679.
- [6] David Cash, Joseph Jaeger, Stanislaw Jarecki, Charanjit S. Jutla, Hugo Krawczyk, Marcel-Catalin Rosu, and Michael Steiner. 2014. Dynamic Searchable Encryption in Very-Large Databases: Data Structures and Implementation. In *NDSS 2014*.
- [7] David Cash, Stanislaw Jarecki, Charanjit S. Jutla, Hugo Krawczyk, Marcel-Catalin Rosu, and Michael Steiner. 2013. Highly-Scalable Searchable Symmetric Encryption with Support for Boolean Queries. In *CRYPTO 2013*. 353–373.
- [8] Javad Ghareh Chamani, Dimitrios Papadopoulos, Charalampos Papamanthou, and Rasool Jalili. 2018. New Constructions for Forward and Backward Private Symmetric Searchable Encryption. In *ACM CCS 2018*. 1038–1055.
- [9] Yan-Cheng Chang and Michael Mitzenmacher. 2005. Privacy Preserving Keyword Searches on Remote Encrypted Data. In *ACNS 2005*. 442–455.
- [10] Melissa Chase and Seny Kamara. 2010. Structured Encryption and Controlled Disclosure. In *ASIACRYPT 2010*. 577–594.
- [11] Reza Curtmola, Juan A. Garay, Seny Kamara, and Rafail Ostrovsky. 2006. Searchable symmetric encryption: improved definitions and efficient constructions. In *ACM CCS 2006*. 79–88.
- [12] Mohammad Etamad, Alptekin Küpçü, Charalampos Papamanthou, and David Evans. 2018. Efficient Dynamic Searchable Encryption with Forward Privacy. *PoPETs 2018*, 1 (2018), 5–20.
- [13] Sky Faber, Stanislaw Jarecki, Hugo Krawczyk, Quan Nguyen, Marcel-Catalin Rosu, and Michael Steiner. 2015. Rich Queries on Encrypted Data: Beyond Exact Matches. In *ESORICS 2015*. 123–145.
- [14] C. Gentry. 2009. Fully homomorphic encryption using ideal lattices. In *ACM STOC'09*. 169–178.
- [15] Eu-Jin Goh. 2003. Secure Indexes. *IACR Cryptology ePrint Archive 2003* (2003), 216.
- [16] Oded Goldreich and Rafail Ostrovsky. 1996. Software Protection and Simulation on Oblivious RAMs. *J. ACM* 43, 3 (1996), 431–473.
- [17] Zichen Gui, Kenneth G. Paterson, Sikhar Patranabis, and Bogdan Warinschi. 2020. SWiSSE: System-Wide Security for Searchable Symmetric Encryption. *IACR Cryptol. ePrint Arch.* (2020), 1328.
- [18] Mohammad Saiful Islam, Mehmet Kuzu, and Murat Kantarcioglu. 2012. Access Pattern disclosure on Searchable Encryption: Ramification, Attack and Mitigation. In *NDSS 2012*.
- [19] Seny Kamara and Tarik Moataz. 2017. Boolean Searchable Symmetric Encryption with Worst-Case Sub-linear Complexity. In *EUROCRYPT 2017*. 94–124.
- [20] Seny Kamara and Tarik Moataz. 2019. Computationally Volume-Hiding Structured Encryption. In *EUROCRYPT 2019*. 183–213.
- [21] Seny Kamara and Charalampos Papamanthou. 2013. Parallel and Dynamic Searchable Symmetric Encryption. In *FC 2013*. 258–274.
- [22] Seny Kamara, Charalampos Papamanthou, and Tom Roeder. 2012. Dynamic searchable symmetric encryption. In *ACM CCS 2012*. 965–976.
- [23] Shangqi Lai, Sikhar Patranabis, Amin Sakzad, Joseph K. Liu, Debdeep Mukhopadhyay, Ron Steinfeld, Shifeng Sun, Dongxi Liu, and Cong Zuo. 2018. Result Pattern Hiding Searchable Encryption for Conjunctive Queries. In *ACM CCS 2018*. 745–762.
- [24] Christopher Manning, Prabhakar Raghavan, and Hinrich Schütze. 2010. Introduction to information retrieval. *Natural Language Engineering* 16, 1 (2010), 100–103.
- [25] Simon Oya and Florian Kerschbaum. 2021. Hiding the Access Pattern is Not Enough: Exploiting Search Pattern Leakage in Searchable Encryption. In *USENIX Security 2021*. 127–142.
- [26] Sarvar Patel, Giuseppe Persiano, Joon Young Seo, and Kevin Yeo. 2021. Efficient Boolean Search over Encrypted Data with Reduced Leakage. In *Advances in Cryptology - ASIACRYPT 2021 (Lecture Notes in Computer Science, Vol. 13092)*, Mehdi Tibouchi and Huaxiong Wang (Eds.). Springer, 577–607. https://doi.org/10.1007/978-3-030-92078-4_20
- [27] Sarvar Patel, Giuseppe Persiano, Joon Young Seo, and Kevin Yeo. 2021. Efficient Boolean Search over Encrypted Data with Reduced Leakage. *IACR Cryptol. ePrint Arch.* (2021), 1227.

- [28] Sarvar Patel, Giuseppe Persiano, Kevin Yeo, and Moti Yung. 2019. Mitigating Leakage in Secure Cloud-Hosted Data Structures: Volume-Hiding for Multi-Maps via Hashing. In *ACM CCS 2019*, Lorenzo Cavallaro, Johannes Kinder, XiaoFeng Wang, and Jonathan Katz (Eds.). ACM, 79–93.
- [29] Sikhar Patranabis and Debdeep Mukhopadhyay. 2021. Forward and Backward Private Conjunctive Searchable Symmetric Encryption. In *NDSS 2021*.
- [30] Dawn Xiaodong Song, David A. Wagner, and Adrian Perrig. 2000. Practical Techniques for Searches on Encrypted Data. In *IEEE S&P 2000*. 44–55.
- [31] Shifeng Sun, Xingliang Yuan, Joseph K. Liu, Ron Steinfeld, Amin Sakzad, Viet Vo, and Surya Nepal. 2018. Practical Backward-Secure Searchable Encryption from Symmetric Puncturable Encryption. In *ACM CCS 2018*. 763–780.
- [32] Yupeng Zhang, Jonathan Katz, and Charalampos Papamanthou. 2016. All Your Queries Are Belong to Us: The Power of File-Injection Attacks on Searchable Encryption. In *USENIX Security Symposium 2016*. 707–720.

A PROOF OF LEMMA 3.1

PROOF OF LEMMA 3.1. We show in the formal proof of Lemma 3.1 that each mkw constructed following the description of Lemma 3.1 covers each w in q ($DB(q)$ part). Other ws that are not in q are filtered out (due to the intersection in Lemma 3.1).

We start with the following conjunctive meta-keyword expression of q_{mkw} for a particular query q as given in the Lemma 3.1.

$$\begin{aligned} DB(q_{mkw}) &= DB(mkw_{i_0, j_0} \wedge \dots \wedge mkw_{i_n, j_n}) \\ &= \bigcap_{k=0}^n DB(mkw_{i_k, j_k}) \end{aligned} \quad (1)$$

By the definition of meta-keywords (Definition 3.1), the following relation holds.

$$DB(mkw_{i_k, j_k}) = \bigcup_{l \in [N] \setminus \{i_k, j_k\}} DB(w_l)$$

We rewrite Equation (1) in the following way.

$$\begin{aligned} DB(q_{mkw}) &= \bigcap_{k=0}^n DB(mkw_{i_k, j_k}) \\ &= \bigcap_{k=0}^n \left(\bigcup_{l \in [N] \setminus \{i_k, j_k\}} DB(w_l) \right) \\ &= \bigcap_{k=0}^n \left(\bigcup_{r \in [n]} DB(w_{\ell_r}) \cup \bigcup_{\substack{l \in [n] \setminus \\ (\{i_k, j_k\}) \\ \cup \{\ell_r : r \in [n]\}}} DB(w_l) \right) \\ &= \bigcup_{r \in [n]} DB(w_{\ell_r}) \cup \bigcap_{k=0}^n \left(\bigcup_{\substack{l \in [N] \\ \setminus (\{i_k, j_k\}) \\ \cup \{\ell_r : r \in [n]\}}} DB(w_l) \right) \\ &= DB(q) \cup \left(\bigcap_{k=0}^n \left(\bigcup_{\substack{l \in [N] \\ \setminus (\{i_k, j_k\}) \\ \cup \{\ell_r : r \in [n]\}}} DB(w_l) \right) \right) \end{aligned}$$

Observe that, the union inside the right hand expression in the above expression keeps all ws except a stretch of ws (from index i_k to j_k) for each value of k , inside the outer intersection of $n + 1$ terms. Since the intersection of these unions reduces to a small but

finite set of id-s, the following relation holds,

$$\mathbf{DB}(q) \subseteq \mathbf{DB}(q_{\text{mkw}})$$

which is exactly what Lemma 3.1 states. \square

B PROOF OF THEOREM 4.1 (CORRECTNESS OF TWINSSE_{BASIC} AND TWINSSE)

The proof of correctness for TWINSSE_{BASIC} (and TWINSSE as well) follows from the correctness of CSSE. The correctness of CSSE ensures that a conjunctive query $q = w_1 \wedge \dots \wedge w_n$ over an encrypted database satisfies the following relations.

$$\mathbf{EDB} = \text{CSSE.SETUP}(\mathbf{DB})$$

$$\mathbf{DB}(w_1) \cap \dots \cap \mathbf{DB}(w_n) = \text{CSSE.SEARCH}(q, \mathbf{EDB})$$

We state the proof for TWINSSE_{BASIC} first. Then we show that this can be simply extended to main TWINSSE scheme (the final bucketized version).

Proof for TWINSSE_{BASIC}. Proof of the TWINSSE_{BASIC} directly follows from the proof of Lemma 3.1. Consider a disjunctive query q as stated below.

$$q = w_1 \vee \dots \vee w_n$$

The equivalent conjunctive expression of meta-keywords can be expressed as below.

$$q_{\text{mkw}} = \text{mkw}_{i_0, j_0} \wedge \text{mkw}_{i_1, j_1} \wedge \dots \wedge \text{mkw}_{i_n, j_n}$$

We write the following relation from Lemma 3.1.

$$\mathbf{DB}(q_{\text{mkw}}) = \mathbf{DB}(q) \cup \left(\bigcap_{k=0}^n \left(\bigcup_{\substack{l \in [N] \setminus \{i_k, j_k\} \\ \cup \{\ell_r : r \in [n]\}}} \mathbf{DB}(w_l) \right) \right)$$

It easy to notice from the above equation that $\mathbf{DB}(q) \subseteq \mathbf{DB}(q_{\text{mkw}})$. Hence, all *ids* of the actual result set of disjunctive query q is included in the result set obtained from the query using TWINSSE_{BASIC}.SEARCH, which proves the correctness of the TWINSSE_{BASIC}.

Proof of TWINSSE. Recall from Section 4, that in TWINSSE all *ws* from Δ are partitioned into n_B buckets of uniform size, and we execute the basic meta-keyword generation method developed in TWINSSE_{BASIC} over each partition independently. Only those partitions with query meta-keywords are accessed during search.

Assume that the dictionary of *ws* - Δ is partitioned in the following way,

$$\Delta = \Delta_1 \cup \Delta_2 \cup \dots \cup \Delta_{n_B}$$

where n_B is the number of buckets and each bucket Δ_u can be expressed in the following way.

$$\Delta_u = \{w_{(u-1)n'+1}, w_{(u-1)n'+2}, \dots, w_{un'}\}$$

The number of *ws* in each bucket is denoted by n' . The set of *mkws* in each bucket Δ_k are represented by the $S_{\text{mkw}, k}$. The TWINSSE_{BASIC} is executed over each of these bucket individually to generate the encrypted database.

The query expression follows from the TWINSSE construction with the above structure (discussed in Section 4).

$$\mathbf{DB}(q_{\text{mkw}}) = \bigcup_{u \in [n_B]} \mathbf{DB}(q_{\text{mkw}, u})$$

and the actual query q can be partitioned in the following way.

$$q = \bigvee_{u \in [n_B]} q_u$$

We expand the above expression to individual buckets.

$$\begin{aligned} \mathbf{DB}(q_{\text{mkw}}) &= \bigcup_{u \in [n_B]} \mathbf{DB}(q_{\text{mkw}, u}) \\ &= \bigcup_{u \in [n_B]} \left(\mathbf{DB}(q_u) \cup \left(\bigcap_{k=0}^{|q_u|} \left(\bigcup_{\substack{l \in [|\Delta_u|] \setminus \\ (\{i_k, j_k\}) \\ \cup \{\ell_r : r \in [q_u]\}}} \mathbf{DB}(w_l) \right) \right) \right) \\ &= \bigcup_{u \in [n_B]} \mathbf{DB}(q_u) \cup \bigcup_{u \in [n_B]} \left(\bigcap_{k=0}^{|q_u|} \left(\bigcup_{\substack{l \in [|\Delta_u|] \setminus \\ (\{i_k, j_k\}) \\ \cup \{\ell_r : r \in [q_u]\}}} \mathbf{DB}(w_l) \right) \right) \\ &= \mathbf{DB}(q) \cup \bigcup_{u \in [n_B]} \left(\bigcap_{k=0}^{|q_u|} \left(\bigcup_{\substack{l \in [|\Delta_u|] \setminus \\ (\{i_k, j_k\}) \\ \cup \{\ell_r : r \in [q_u]\}}} \mathbf{DB}(w_l) \right) \right) \end{aligned}$$

Clearly, from the above expression $\mathbf{DB}(q) \subseteq \mathbf{DB}(q_{\text{mkw}})$, where $|\mathbf{DB}(q)| = \eta \cdot |\mathbf{DB}(q_{\text{mkw}})|$ (recall from Section 3.3 that η is the precision parameter), which proves the correctness of result for TWINSSE.

C ADAPTIVE SECURITY OF SSE

The adaptive security of any SSE scheme is parameterized by a leakage function

$$\mathcal{L} = \left(\mathcal{L}^{\text{SETUP}}, \mathcal{L}^{\text{SEARCH}} \right),$$

where $\mathcal{L}^{\text{SETUP}}$ encapsulates the leakage to an adversarial server during the setup phase, and $\mathcal{L}^{\text{SEARCH}}$ encapsulates the leakage to an adversarial server during each execution of the search protocol.

Algorithm 5 Experiment $\text{Real}^{\text{SSE}}(\lambda, Q)$

```

1: function  $\text{Real}^{\text{SSE}}(\lambda, Q)$ 
2:    $N \leftarrow \text{Adv}(\lambda)$ 
3:    $(\text{sk}, \text{st}_0, \mathbf{EDB}_0) \leftarrow \text{SETUP}(\lambda, N)$ 
4:   for  $k \leftarrow 1$  to  $Q$  do
5:     Let  $q_k \leftarrow \text{Adv}(\lambda, \mathbf{EDB}_{k-1}, \tau_1, \dots, \tau_{k-1})$ 
6:     Let  $(\text{st}_k, \mathbf{EDB}_k, \mathbf{DB}(q_k)) \leftarrow$ 
       SEARCH( $\text{sk}, \text{st}_{k-1}, q_k; \mathbf{EDB}_{k-1}$ )
7:     Let  $\tau_k$  denote the view of the adversary after
       the  $k^{\text{th}}$  query
8:    $b \leftarrow \text{Adv}(\lambda, \mathbf{EDB}_Q, \tau_1, \dots, \tau_Q)$ 
9:   return  $b$ 

```

Informally, an SSE scheme is adaptively secure with respect to a leakage function \mathcal{L} if the adversarial server provably learns no more information about \mathbf{DB} other than that encapsulated by \mathcal{L} . Formally,

Algorithm 6 Experiment $\mathbf{Ideal}^{\text{SSE}}(\lambda, Q, \mathcal{L})$

```

1: function  $\mathbf{Ideal}^{\text{SSE}}(\lambda, Q, \mathcal{L})$ 
2:   Parse the leakage function  $\mathcal{L}$  as:
      $\mathcal{L} = (\mathcal{L}^{\text{SETUP}}, \mathcal{L}^{\text{SEARCH}})$ .
3:    $(\text{st}_{\text{SIM}}, \mathbf{EDB}_0) \leftarrow \text{SIMSETUP}(\mathcal{L}^{\text{SETUP}}(\lambda, N))$ 
4:   for  $k \leftarrow 1$  to  $Q$  do
5:     Let  $q_k \leftarrow \text{Adv}(\lambda, \mathbf{EDB}_{k-1}, \tau_1, \dots, \tau_{k-1})$ 
6:     Let  $(\text{st}_{\text{SIM}}, \mathbf{EDB}_k, \tau_k) \leftarrow \text{SIMSEARCH}$ 
        $(\text{st}_{\text{SIM}}, \mathcal{L}^{\text{SEARCH}}(q_k); \mathbf{EDB}_{k-1})$ 
7:     Let  $\tau_k$  denote the view of the adversary after
       the  $k^{\text{th}}$  query
8:    $b \leftarrow \text{Adv}(\lambda, \mathbf{EDB}_Q, \tau_1, \dots, \tau_Q)$ 
9:   return  $b$ 

```

an SSE scheme is said to be adaptively secure with respect to a leakage function \mathcal{L} if for any stateful PPT adversary Adv that issues a maximum of $Q = \text{poly}(\lambda)$ queries, there exists a stateful probabilistic polynomial-time simulator $\text{SIM} = (\text{SIMSETUP}, \text{SIMSEARCH})$ such that the following holds:

$$\left| \Pr \left[\mathbf{Real}_{\text{Adv}}^{\text{SSE}}(\lambda, Q) = 1 \right] - \Pr \left[\mathbf{Ideal}_{\text{Adv, SIM}}^{\text{SSE}}(\lambda, Q) = 1 \right] \right| \leq \text{negl}(\lambda),$$

where the “real” experiment $\mathbf{Real}^{\text{SSE}}$ and the “ideal” experiment $\mathbf{Ideal}^{\text{SSE}}$ are as described in Algorithm 5 and Algorithm 6.

D DETAILED ANALYSIS AND DISCUSSION ON THE LEAKAGE OF TWINSSE_{OXT}

In Section 4.4, we informally described the leakage profile for TWINSSE built in a black-box way from any generic conjunctive SSE scheme CSSE. In this section, we formally detail the leakage profile for the specific instantiation of TWINSSE based on the OXT scheme, namely TWINSSE_{OXT}. We then present a discussion on the leakage profiles for both TWINSSE and TWINSSE_{OXT}.

D.1 Security of TWINSSE

We present a formal description of the security guarantees of our generic construction TWINSSE. Concretely, we state the following theorem.

THEOREM D.1 (SECURITY OF TWINSSE). *Assuming that CSSE is an (adaptively) secure SSE scheme with respect to the leakage function $\mathcal{L}_{\text{CSSE}} = (\mathcal{L}_{\text{CSSE}}^{\text{SETUP}}, \mathcal{L}_{\text{CSSE}}^{\text{SEARCH}})$, TWINSSE is an (adaptively) secure SSE scheme with respect to the leakage function $\mathcal{L}_{\text{TWINSSE}} = (\mathcal{L}_{\text{TWINSSE}}^{\text{SETUP}}, \mathcal{L}_{\text{TWINSSE}}^{\text{SEARCH}})$, where for any plaintext database \mathbf{DB} , any search query q , and any pair of bucketization parameters (n', n_B) , we have*

$$\mathcal{L}_{\text{TWINSSE}}^{\text{SETUP}}(\mathbf{DB}) = \left(\mathcal{L}_{\text{CSSE}}^{\text{SETUP}}(\widehat{\mathbf{DB}}), n', n_B \right),$$

where $\widehat{\mathbf{DB}} = \text{GENMETADB}(\mathbf{DB}, n', n_B)$, and

$$\mathcal{L}_{\text{TWINSSE}}^{\text{SEARCH}}(q) = \begin{cases} \mathcal{L}_{\text{CSSE}}^{\text{SEARCH}}(q) & \text{if } q \text{ is conjunctive,} \\ \left\{ \mathcal{L}_{\text{CSSE}}^{\text{SEARCH}}(q_{\text{mkw}, k}) \right\}_{k \in [n_B]} & \text{if } q \text{ is disjunctive,} \end{cases}$$

where

$$q_{\text{mkw}} = \bigvee_{k \in [n_B]} q_{\text{mkw}, k} = \text{GENMQUERY}(q, n', n_B).$$

Proof. We defer the formal proof of this theorem to Appendix E.

D.2 Leakage Profile of TWINSSE_{OXT}

In this section, we describe the leakage profile of TWINSSE_{OXT}. We begin by recalling from [7] the leakage profile of the original OXT scheme. We then build upon it to describe the leakage profile of TWINSSE_{OXT}, which is actually very similar in spirit to the leakage profile of OXT.

Setup Leakage. The setup leakage in the OXT scheme consists of the size of the database \mathbf{DB} , which is nothing but the total number of keyword-document pairs in the database \mathbf{DB} , formally defined as

$$|\mathbf{DB}| = \sum_{w \in \Delta} |\mathbf{DB}(w)|,$$

where $\Delta = \{w_1, \dots, w_N\}$ is the dictionary over which the database \mathbf{DB} is defined.

Search Leakages. Next, we summarize the leakages incurred by OXT during conjunctive keyword search queries.

Result Pattern Leakage: The server learns the final set of document identifiers matching the query. Formally, for a conjunctive query q , the result pattern leakage RP is defined as

$$\text{RP}(q) = \mathbf{DB}(q).$$

Size Pattern Leakage. The server learns the frequency of the s -term (where s -term again refers to the least frequent keyword in the conjunction). Formally, for a conjunctive query $q = (w_1 \wedge \dots \wedge w_n)$, where w_1 is the least frequent keyword in the conjunction, the size pattern SP is defined as

$$\text{SP}(q) = |\mathbf{DB}(w_1)|.$$

Equality Pattern Leakage. The server learns if two (or more) conjunctive queries have the same s -term (where s -term again refers to the least frequent keyword in the conjunction). Formally, for a sequence of conjunctive queries (q_1, \dots, q_M) , where for each $i \in [M]$, we have

$$q_i = (w_{i,1} \wedge \dots \wedge w_{i,n_i}),$$

where $w_{i,1}$ is the least frequent keyword in the conjunction, the equality pattern leakage EP is defined as an $M \times M$ matrix where for each $i, j \in [M]$, we have

$$\text{EP}[i, j] = \begin{cases} 1 & \text{if } w_{i,1} = w_{j,1}, \\ 0 & \text{otherwise.} \end{cases}$$

Conditional Intersection Pattern Leakage. The server learns if two (or more) conjunctive queries have one or more x -terms in common (where x -term refers any keyword other than the least frequent keyword in the conjunction); more concretely, if two (or more) conjunctive queries have one or more x -terms in common, then the server learns the intersection of the set of documents containing the corresponding s -terms. Formally, for a sequence of conjunctive queries (q_1, \dots, q_M) , where for each $i \in [M]$, we have

$$q_i = (w_{i,1} \wedge \dots \wedge w_{i,n_i}),$$

where $w_{i,1}$ is the least frequent keyword in the conjunction, the conditional intersection pattern leakage IP is defined as an $M \times M$ matrix of lists, where for each $i, j \in [M]$, we have

$$\text{IP}[i, j] = \begin{cases} \text{DB}(w_{i,1}) \cap \text{DB}(w_{j,1}) & \text{if } \overline{\text{IP}}[i, j] = 1, \\ \phi & \text{if } \overline{\text{IP}}[i, j] = 0, \end{cases}$$

where $\overline{\text{IP}}[i, j] = 1$ if and only if there exists at least one pair $(\ell_i, \ell_j) \in [n_i] \times [n_j]$ such that $w_{i,\ell_i} = w_{j,\ell_j}$; otherwise, we have $\overline{\text{IP}}[i, j] = 0$.

Security of TWINSSE_{OXT}. We now formalize the security of TWINSSE_{OXT} in terms of the leakage profiles described above. We do this using a formal theorem, which may be viewed as a specialization of Theorem D.1 to a specific instantiation of TWINSSE based on OXT. Once again, this theorem is based on the (adaptive) simulation-security definition of SSE in the real world-ideal world paradigm, described in Appendix C.

THEOREM D.2 (SECURITY OF TWINSSE_{OXT}). TWINSSE_{OXT} is an (adaptively) secure SSE scheme with respect to the leakage function $\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}} = (\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SETUP}}, \mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH}})$, where for any plaintext database DB , any sequence of conjunctive queries $Q_0 = (q_{1,0}, \dots, q_{M,0})$ and any sequence of disjunctive queries $Q_1 = (q_{1,1}, \dots, q_{M',1})$, and any pair of bucketization parameters (n', n_B) , we have

$$\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SETUP}}(\text{DB}) = (\widehat{\text{DB}}, n', n_B),$$

where $\widehat{\text{DB}} = \text{GENMETADB}(\text{DB}, n', n_B)$, and

$$\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH}}(Q_0, Q_1) = [\text{RP}, \text{SP}, \text{EP}, \text{IP}](Q_0, Q_{\text{mkw},1}),$$

where $Q_{\text{mkw},1}$ is a sequence of (sub-)meta-queries of the form

$$Q_{\text{mkw},1} = \{q_{\text{mkw},k,\ell}\}_{k \in [n_B], \ell \in [M']},$$

where for each $\ell \in [M']$, we have

$$q_{\text{mkw},\ell} = \left(\bigvee_{k \in [n_B]} q_{\text{mkw},k,\ell} \right) = \text{GENMQUERY}(q_{\ell,1}, n', n_B).$$

Proof. We defer the formal proof of this theorem to Appendix E.

D.3 Discussion on the Leakage Profile of TWINSSE_{OXT}

In this subsection, we present a more in-depth analysis of the leakage profile for TWINSSE_{OXT} during conjunctive and disjunctive search queries, and its implications.

Output Leakage. We begin by noting that the output leakage (alternatively, the result pattern leakage) is incurred by nearly all existing SSE schemes, including static and dynamic schemes, in the setting of both single and conjunctive keyword searches (such as in [4, 7, 8, 11, 23, 31]). This is usually considered acceptable in the SSE literature; indeed the few known data/query recovery attacks that manage to exploit this leakage ([2, 5, 18, 32]) assume extremely strong adversarial models where the adversary has partial knowledge of the plaintext database/queries.

s-Term Leakages. We focus next on the leakages related to the s -term, namely the size and equality pattern leakages. We begin by noting that these leakages are somewhat inherent to the design paradigm of OXT, which we base our instantiation of TWINSSE on. Even in the simpler setting of single keyword search, most existing

schemes [3, 4, 6, 8, 10, 11, 31] also incur size and equality pattern leakages; the only constructions not to incur such leakages seem to rely on the use of ORAM-style data structures [4, 8]. Fortifying TWINSSE_{OXT} with such data structures in an attempt to prevent this leakage is an interesting open challenge, although this would probably have to trade-off with some degradation in search performance (mostly in terms of communication complexity and number of rounds of communication during searches).

It is also possible (and perhaps conceptually simpler) to mask this leakage by using volume-hiding techniques such as padding and encrypted multi-maps (EMMs) [1, 11, 20, 27, 28]. This would incur a degradation in search performance, and it is up to the designer to decide on a suitable trade-off between performance and leakage.

However, we would like to point out that there are no known data/query recovery attacks on SSE schemes that specially exploit leakages related to the s -term. So we believe that *even without the aforementioned fortifications*, it appears that our TWINSSE_{OXT} scheme is not vulnerable to any known attacks due to the leakages related to the s -term, in realistic/practical adversarial settings.

x-Term Leakages. Next, we focus on the x -term leakages. We again note that these leakages are somewhat inherent to the design paradigm of OXT, which we base our instantiation of TWINSSE on. The only known attack on conjunctive SSE schemes that exploits a form of x -term leakages is the *file injection attack* proposed by Zhang *et al.* in [32]. More concretely, the adversarial server must be able to compute $|\text{DB}(w_1) \cap \text{DB}(w_i)|$ when processing the search query.

We note however that for file injection attacks to work efficiently, the adversarial server must recover, for every x -term w_i , the result size corresponding to each sub-query of the form $w_1 \cap w_i$. However, the x -term leakage profile of TWINSSE_{OXT} is not sufficient to compute this term, since the set of x token values sent to the server is randomly permuted inside the underlying OXT instantiation precisely to mask such inference-style attacks. Further, one could also instantiate our generic construction of TWINSSE from other conjunctive SSE schemes such as HXT [23] that improve upon OXT in terms of provable security against leakage-based cryptanalysis and file-injection attacks.

Finally, fortifying implementations of TWINSSE_{OXT} by using either ORAM-style data structures or adopting volume-hiding techniques such as padding/EMMs may be useful in masking this leakage even further; however, even without such additional fortifications, it appears that our TWINSSE_{OXT} scheme is not vulnerable to file injection attacks, or any other known attacks for that matter, due to the leakages related to the s -term, in realistic/practical adversarial settings.

Leakage Cryptanalysis. Looking ahead, in Appendix F, we present a leakage-based cryptanalysis of the TWINSSE_{OXT} scheme via experiments over the Enron email corpus. Our experiments help substantiate that the leakages incurred by the disjunctive search protocol in TWINSSE_{OXT} are reasonably benign in practice and are quite resistant to even the most powerful leakage-based cryptanalysis techniques in the SSE literature over real-world databases, such as those in [5, 32]. We leave it as an open question to extend the analysis using the more advanced leakage cryptanalysis techniques, such as those proposed in [2, 25].

E SECURITY PROOFS FOR TWINSSE AND TWINSSE_{OXT}

In this section, we formally prove the security of TWINSSE and TWINSSE_{OXT} with respect to the generic and specific leakage profiles described in Theorems D.1 and D.2, respectively.

E.1 Proof of Theorem D.1 (Security Analysis of TWINSSE)

We provide a simulation-based proof approach for TWINSSE. We assumed that the underlying adaptively secure CSSE has the following leakage profile.

$$\mathcal{L}_{\text{CSSE}} = (\mathcal{L}_{\text{CSSE}}^{\text{SETUP}}, \mathcal{L}_{\text{CSSE}}^{\text{SEARCH}})$$

We express the leakage of TWINSSE as,

$$\mathcal{L}_{\text{TWINSSE}} = (\mathcal{L}_{\text{TWINSSE}}^{\text{SETUP}}, \mathcal{L}_{\text{TWINSSE}}^{\text{SEARCH}})$$

where,

$$\mathcal{L}_{\text{TWINSSE}}^{\text{SETUP}}(\mathbf{DB}) = \mathcal{L}_{\text{CSSE}}^{\text{SETUP}}(\widehat{\mathbf{DB}})$$

and, $\widehat{\mathbf{DB}} = \text{GENMETADB}(\mathbf{DB}, n', n_B)$, and

$$\mathcal{L}_{\text{TWINSSE}}^{\text{SEARCH}}(q) = \begin{cases} \mathcal{L}_{\text{CSSE}}^{\text{SEARCH}}(q) & \text{if } q \text{ is conjunctive,} \\ \{\mathcal{L}_{\text{CSSE}}^{\text{SEARCH}}(q_{\text{mkw},k})\}_{k \in [n_B]} & \text{if } q \text{ is disjunctive,} \end{cases}$$

where

$$q_{\text{mkw}} = \left(\bigvee_{k \in [n_B]} q_{\text{mkw},k} \right) = \text{GENMQUERY}(q, n', n_B).$$

We show that TWINSSE is secure against an adaptive semi-honest adversary \mathcal{A} , which has access to leakages from TWINSSE. We build a simulator $\text{SIM}_{\widehat{\mathbf{EDB}}}$ generation by TWINSSE.SETUP , and transcripts for queries over $\widehat{\mathbf{EDB}}$. The simulator simulates the transcripts τ_i for each query q_i . The simulator has the inputs from the leakage function $\mathcal{L}_{\text{TWINSSE}}$ only, with the setup leakage $\mathcal{L}_{\text{TWINSSE}}^{\text{SETUP}}$ and the search leakage $\mathcal{L}_{\text{TWINSSE}}^{\text{SEARCH}}$.

Simulating TWINSSE.SETUP: The following public parameters are available to SIM_{CSSE} as a part of $\mathcal{L}_{\text{TWINSSE}}^{\text{SETUP}}$.

$$\{\mathbf{DB}, n', n_B\}$$

The simulator outputs the its version of $\widehat{\mathbf{EDB}}$ according to the simulation process of CSSE (we assumed that CSSE is provably simulation secure).

$$\begin{aligned} ct_{\widehat{\mathbf{EDB}}} &= \text{SIM}_{\text{TWINSSE}}^{\text{SETUP}}(\mathbf{DB}) \\ &= \text{SIM}_{\text{CSSE}}^{\text{SETUP}}(\widehat{\mathbf{DB}}) \\ &= \text{SIM}_{\text{CSSE}}^{\text{SETUP}}(\mathbf{DB}, n', n_B) \end{aligned}$$

Since, CSSE is proven simulation secure, it follows from the simulation security guarantee of CSSE that $ct_{\widehat{\mathbf{EDB}}}$ is indistinguishable from the one generated in the real experiment.

Simulating TWINSSE.SEARCH: For conjunctive queries the adversary does not have any advantage from $\mathcal{L}_{\text{TWINSSE}}^{\text{SEARCH}}$ compared to $\mathcal{L}_{\text{CSSE}}^{\text{SEARCH}}$, which exactly same as CSSE. For disjunctive queries we consider the effect of querying using q_{mkw} .

For disjunctive queries, we argue that the adversary \mathcal{A} does not gain any information about the original disjunctive query with

this simulation experiment. The distribution of $\widehat{\mathbf{DB}}$ (hence, also for $\widehat{\mathbf{EDB}}$) is abstracted from \mathbf{DB} by the meta-keywords. The search leakages of CSSE is characterised by the $\mathcal{L}_{\text{CSSE}}$, provided from CSSE construction. Since, CSSE in TWINSSE executes over meta-keyword only, this leakage is expressed in the context of meta-keywords as below.

$$\mathcal{L}'_{\text{CSSE}} = \mathcal{L}_{\text{CSSE}}(\text{meta-keywords})$$

With this leakage information of CSSE, the search leakage of TWINSSE can be expressed as below.

$$\mathcal{L}_{\text{TWINSSE}}^{\text{SEARCH}}(q) = \mathcal{L}_{\text{TWINSSE}}^{\text{SEARCH}}(q_{\text{mkw},k})_{k \in [n_B]} = \{\mathcal{L}'_{\text{CSSE}}, n_B, n'\}$$

The parameters n_B and n' are derived from N (number of keywords), which is available during setup. Therefore, the search leakage of TWINSSE same as the underlying CSSE, which can be summarised as below.

$$\mathcal{L}_{\text{TWINSSE}}^{\text{SEARCH}}(q) = \mathcal{L}_{\text{TWINSSE}}^{\text{SEARCH}}(q_{\text{mkw},k})_{k \in [n_B]} = \{\mathcal{L}'_{\text{CSSE}}\}$$

This same leakage profile for search in TWINSSE and CSSE in the context of meta-keywords ensures that no additional information is leaked beyond CSSE leakage.

E.2 Proof of Theorem D.2 (Security Analysis of TWINSSE_{OXT})

We resort to a simulation-based security analysis for TWINSSE_{OXT}. We assume a semi-honest adversary \mathcal{A} which has access to the leakage from standard SSE leakages in an adaptive model. Security analysis of TWINSSE relies upon the semantic security notions provided by CSSE. TWINSSE inherits these notions through the core OXT (in case of TWINSSE_{OXT}, the OXT) instance. We assume the following properties of OXT achieves with efficient performance.

- (1) Primitives used in construction of OXT hold the standard security assumptions.
- (2) OXT is non-adaptively and adaptively secure with the above assumptions.

We consider the following leakage profile for OXT.

$$\mathcal{L}_{\text{OXT}} = \{\mathcal{L}_{\text{OXT}}^{\text{SETUP}}, \mathcal{L}_{\text{OXT}}^{\text{SEARCH}}\}$$

Here, $\mathcal{L}_{\text{OXT}}^{\text{SETUP}}$ captures the leakage from the OXT.SETUP, and $\mathcal{L}_{\text{OXT}}^{\text{SEARCH}}$ encapsulates the leakage from OXT.SEARCH. More precisely, these can be expressed as,

$$\mathcal{L}_{\text{OXT}}^{\text{SETUP}}(\mathbf{DB}) = \{|\mathbf{DB}|\}$$

and

$$\mathcal{L}_{\text{OXT}}^{\text{SEARCH}}(\mathbf{EDB}, \{q_k\}_{q_k \in \mathcal{Q}_0}) = \{RP, SP, EP, IP\}$$

where, \mathcal{Q}_0 is a set of conjunctive queries. The leakages RP , SP , EP , and IP are the pattern leakages from OXT (see Appendix D.3).

We define the leakage profile of TWINSSE_{OXT} with respect to these above definitions and assumptions as below.

$$\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}} = \{\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SETUP}}, \mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH}}\}$$

The leakage functions above can be expressed as

$$\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SETUP}}(\mathbf{DB}) = \{|\widehat{\mathbf{DB}}|, n', n_B\}$$

and

$$\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH}}(\widehat{\text{EDB}}, Q_0, Q_1) = [\text{RP}, \text{SP}, \text{EP}, \text{IP}](Q_0, Q_{\text{mkw},1}).$$

For conjunctive queries,

$$\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH}}(\widehat{\text{EDB}}, \{q_k\}_{q_k \in Q_0}) = [\widehat{\text{RP}}, \widehat{\text{SP}}, \widehat{\text{EP}}, \widehat{\text{IP}}].$$

Here, $\{\widehat{\text{RP}}, \widehat{\text{SP}}, \widehat{\text{EP}}, \widehat{\text{IP}}\}$ are the $\{\text{RP}, \text{SP}, \text{EP}, \text{IP}\}$ leakages in the context of meta-keywords. For conjunctive queries, it is exactly the same as OXT.

Since, OXT is simulation secure against these leakages, simulation security of $\text{TWINSSE}_{\text{OXT}}$ for conjunctive queries is straightforwardly implied from OXT.

In disjunctive queries, the query transformation process is carried out locally by the client, and the actual search is completed using OXT.SEARCH protocol, we can write $\text{TWINSSE}_{\text{OXT}.SEARCH}$ leakage as

$$\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH}}(\widehat{\text{EDB}}, \{q_{\text{mkw},1,k}\}_{k \in [Q_1]}) = \{\widehat{\text{RP}}, \widehat{\text{SP}}, \widehat{\text{EP}}, \widehat{\text{IP}}\}.$$

We build a simulator SIM to simulate the $\widehat{\text{EDB}}$ generation by $\text{TWINSSE}_{\text{OXT}}$ from DB , and transcripts for query search over $\widehat{\text{EDB}}$. The simulator simulates the transcripts τ_i for each query $q_i \in Q$. The simulator has the inputs from the leakage function $\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}$ only, with the setup leakage $\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SETUP}}$ and the search leakage $\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH}}$.

Simulating Setup: The following public parameters are available to SIM_{OXT} as a part of $\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SETUP}}$.

$$\{|\text{EDB}|, |\widehat{\Delta}|\}$$

The simulator outputs the its version of $\widehat{\text{EDB}}$ according to the simulation process of OXT (we assumed that OXT is provably simulation secure).

$$ct_{\widehat{\text{EDB}}} = SIM_{\text{OXT}}.SETUP(|\text{MDB}|, |\widehat{\Delta}|)$$

It follows from the simulation security guarantee of OXT that $ct_{\widehat{\text{EDB}}}$ is indistinguishable from the one generated in the real experiment.

Simulating Search: For the conjunctive queries, the leakage $\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH}}$ is exactly the same as $\mathcal{L}_{\text{OXT}}^{\text{SEARCH}}$. Hence, we can write the following.

$$\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH}}(\widehat{\text{EDB}}, \{q_k\}_{k \in [Q]}) = \mathcal{L}_{\text{OXT}}^{\text{SEARCH}}(\text{EDB}, \{q_k\}_{k \in [Q]})$$

By the simulation security guarantee of OXT, $\text{TWINSSE}_{\text{OXT}}$ secure against these leakages.

For disjunctive queries, we argue that the adversary \mathcal{A} does not gain any information about the original disjunctive query except $|q|$. The distribution of MDB (encrypted to $\widehat{\text{EDB}}$) is abstracted from DB through the meta-keywords. We resort to a more conservative analysis for this proof, as keywords do not have direct inference from meta-keywords, especially that is applicable over any database in general. The position of each w in an mkw is fixed according to the frequency of w , which is unique for a DB . The lemmas below relate worst cases where an inference can be established between the query keywords and the corresponding meta-keywords without any additional knowledge of the plain database.

Lemma E.1, Lemma E.2, and Lemma E.3 relates the disjunctive q with $w_i \in \Delta$ to the conjunctive q with $\text{mkw}_i \in \widehat{\Delta}$.

LEMMA E.1. Consider two disjunctive queries of the same length t

$$q_0 = w_{1,q_0} \vee w_{2,q_0} \vee \dots \vee w_{t,q_0}, w_{i,q_0}$$

$$q_1 = w_{1,q_1} \vee w_{2,q_1} \vee \dots \vee w_{t,q_1}, w_{i,q_1}$$

have the following expressions using mkws ,

$$q_0 = q_{0,\text{mkw}} = \text{mkw}_{1,q_0} \wedge \text{mkw}_{2,q_0} \wedge \dots \wedge \text{mkw}_{t+1,q_0}$$

$$q_1 = q_{1,\text{mkw}} = \text{mkw}_{1,q_1} \wedge \text{mkw}_{2,q_1} \wedge \dots \wedge \text{mkw}_{t+1,q_1}$$

both of length $t + 1$, and the mkws are placed in the increasing order of the starting index of the 0s stretch in each mkw . If the mkws at index k in q_0 and q_1 are the same, then $w_{k-1,q_0} = w_{k-1,q_1}$ and $w_{k,q_0} = w_{k,q_1}$.

PROOF. The proof of Lemma E.1 is given in Section E.3.1. \square

LEMMA E.2. Consider two disjunctive queries q_0 and q_1 , of the same length t have the mkw expressions as defined in Lemma E.1 - both of length $t + 1$. If the mkws at indices k_0 in q_0 , and k_1 in q_1 are the same, then $w_{k_0-1,q_0} = w_{k_1-1,q_1}$ and $w_{k_0,q_0} = w_{k_1,q_1}$.

PROOF. The proof of Lemma E.2 is given in Section E.3.2. \square

LEMMA E.3. Consider two disjunctive queries of different length t_0 and t_1 -

$$q_0 = w_{1,q_0} \vee w_{2,q_0} \vee \dots \vee w_{t_0,q_0}, w_{i,q_0} \in \Delta$$

$$q_1 = w_{1,q_1} \vee w_{2,q_1} \vee \dots \vee w_{t_1,q_1}, w_{i,q_1} \in \Delta$$

have following expressions in the mkws

$$q_0 = q_{0,\text{mkw}} = \text{mkw}_{1,q_0} \wedge \text{mkw}_{2,q_0} \wedge \dots \wedge \text{mkw}_{t_0+1,q_0}$$

$$q_1 = q_{1,\text{mkw}} = \text{mkw}_{1,q_1} \wedge \text{mkw}_{2,q_1} \wedge \dots \wedge \text{mkw}_{t_1+1,q_1}$$

which are of lengths $t_0 + 1$ and $t_1 + 1$ respectively. If the mkws at indices k_0 in q_0 , and k_1 in q_1 are the same, then $w_{k_0-1,q_0} = w_{k_1-1,q_1}$ and $w_{k_0,q_0} = w_{k_1,q_1}$.

PROOF. The proof of Lemma E.3 is given in Section E.3.3. \square

Recall that, the query transformation is executed by the client locally. The search is executed as a two-party protocol between the client and the server using the meta-keywords. The server learns $|q|$ trivially from q_{mkw} through of meta-keywords. From Lemma E.1, E.2, and E.3, an adversary can infer the position of the same w s in two queries of same length or different lengths if both queries have a common mkw in them.

However, the server can only infer if the least-frequent mkws in q_{mkw} are identical or not in mkw expressions of two qs from \widehat{SP} . The mkw expressions in each of the three lemmas require to place mkws in increasing order of the starting index of the 0's stretch. Whereas, the actual query expression for OXT has the least-frequent mkw first. No direct inference can be conjectured for the least-frequent mkw and the query expressions in the lemmas. Hence, an adversary \mathcal{A} can not distinguish between the common meta-keyword and a distinct meta-keyword.

In the case, where the least-frequent of mkws is the first one in the query expression of the lemmas too, the first keyword is also the same for both w s. This is equivalent to the case of two conjunctive queries in keywords having the least-frequent w same.

Therefore, the leakage from $TWINSSE_{SEARCH}$ can be limited to the OXT pattern leakages only, as expressed below.

$$\mathcal{L}_{TWINSSE_{OXT}}^{SEARCH}(\widehat{EDB}, \{q_k\}_{k \in [|Q|]}) = \{\mathcal{L}_{OXT}^*, |q_k|_{k \in [|Q|]}\}$$

Since, OXT is proven simulation secure, it follows from the simulation security guarantee that \mathcal{A} no additional advantage over the real experiment.

E.3 Proofs of the Lemmas

We present the proofs of the lemmas presented earlier in this section. We follow the notations and conventions as used in the main body of the paper.

E.3.1 Proof of Lemma E.1.

PROOF. By construction, each meta-keyword mkw_i has the original keywords appearing in sorted order in the binary string representation (increasing order of frequency from left to right). Assume, the k 'th meta-keyword mkw_k is same for both the queries q_0 and q_1 . Without loss of generality, a meta-keyword in the basic $O(N^2)$ ($TWINSSE_{BASIC}$) method can be formed as

$$\{b_1, b_2, \dots, b_r, b_{r+1}, \dots, b_s, b_{s+1}, \dots, b_n\}, b_i \in \{0, 1\}$$

where $1 \leq r < s \leq n$, and $b_i = 0$ for $r < i < s$.

To have an mkw of this form, q must have two keywords at indices r and s , and none in between (for q_0 and q_1 both). Since the $mkws$ are constructed using ws in sorted order, if both queries q_0 and q_1 have the same r and same s (as one mkw is the same), the keywords w_r and w_s in both q_0 and q_1 are also the same. Hence, we have $w_{k-1, q_0} = w_{k-1, q_1}$ and $w_{k, q_0} = w_{k, q_1}$. \square

E.3.2 Proof of Lemma E.2.

PROOF. We assume the common mkw of q_0 and q_1 can be expressed as

$$\{b_1, b_2, \dots, b_r, b_{r+1}, \dots, b_s, b_{s+1}, \dots, b_n\}, b_i \in \{0, 1\}$$

where $1 \leq r < s \leq n$, and $b_i = 0$ for $r < i < s$. The mkw appears at indices k_0 in q_0 and at k_1 in q_1 . Since the indices of ws in the mkw strings are in sorted order (increasing frequency) and remains fixed for all $mkws$, the ws at index r and index s are the same for both q_0 and q_1 . However, as the index of mkw is different in q_0 and q_1 , the number of preceding ws before index r in q_0 and q_1 are different, equal to $k_0 - 2$ and $k_1 - 2$ respectively. Hence, for q_0 , r is equal to $k_0 - 1$, and equal to $k_1 - 1$ in q_1 . Following the above argument, we have $w_{k_0-1, q_0} = w_{k_1-1, q_1}$ and $w_{k_0, q_0} = w_{k_1, q_1}$. \square

E.3.3 Proof of Lemma E.3.

PROOF. The proof of Lemma E.3 follows from the proof of Lemma E.2. Essentially, Lemma E.3 is the extension of Lemma E.2 for two different lengths of queries. Intuitively, it can be established in the following way. Recall that in Proof E.3.2, r and s remains same in both q_0 and q_1 , as in binary representation all $mkws$ and qs have the same length n . However, the number of ws in q changes, and consequently, number of $mkws$ change. Hence, the range of indices k_0 and k_1 are different for q_0 and q_1 . This does not affect r and s which are positions of keywords (not related to number of keywords) in the binary representation of fixed length. Hence, the same argument from the proof of Lemma E.2 holds. \square

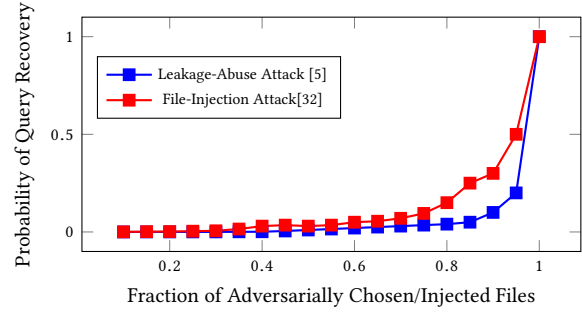


Figure 11: Leakage Analysis of $TWINSSE_{OXT}$: Two-Keyword Conjunctive Searches in the “Chosen Files” Setting.

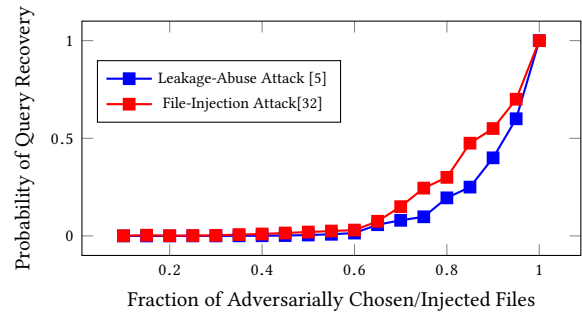


Figure 12: Leakage Analysis of $TWINSSE_{OXT}$: Two-Keyword Disjunctive Searches in the “Chosen Files” Setting.

F CRYPTANALYSIS OF $TWINSSE_{OXT}$

In this section, we present a leakage-based cryptanalysis of the $TWINSSE_{OXT}$ scheme via experiments over the Enron email corpus. Our experiments help substantiate that the leakages incurred by the conjunctive and disjunctive search protocols in $TWINSSE_{OXT}$ are benign in practice and are resistant to even the most powerful leakage-based cryptanalysis techniques in the SSE literature over real-world databases, such as *leakage-abuse attacks* [5] and *file-injection attacks* [32]. In particular, we experimentally establish the following claim:

CLAIM F.1 (INFORMAL). *The conjunctive and disjunctive search protocols in $TWINSSE_{OXT}$ resist leakage-abuse attacks [5] and file-injection attacks [32], and are benign in practice.*

We substantiate this claim by leakage cryptanalysis experiments targeting the conjunctive and disjunctive keyword search protocols in $TWINSSE_{OXT}$. We evaluate the probability that the adversary guesses correctly the keywords w_1 and w_2 underlying a two-conjunction query $q = (w_1 \wedge w_2)$ (resp., a two-disjunction query $q = (w_1 \vee w_2)$) by one of two well-known and extensively studied cryptanalysis methodologies in the SSE literature- the leakage-abuse attack of Cash *et al.* [5] and the *file-injection attack* of Zhang *et al.* [32]. The experiments were conducted over the same Enron email corpus as was used for the performance evaluation experiments in Section 5. The attacks operate in the *chosen/injected* file

model (the strongest possible attack setting where a certain fraction of the files in the database are adversarially generated.) The corresponding results are plotted in Figure 11 and Figure 12 for conjunctive and disjunctive queries, respectively. Throughout, we use a bucket size $n' = 10$ (same as for the performance evaluation experiments in Section 5) for the disjunctive experiments.

We note here that fortifying implementations of $TWINSSE_{OXT}$ by using either ORAM-style data structures or adopting volume-hiding techniques such as padding/encrypted multi-maps [20, 28] may be useful in masking leakage even further. However, even without such additional fortifications, $TWINSSE_{OXT}$ resists leakage-abuse and file-injection attacks in the strongest possible attacker setting, as demonstrated by the aforementioned experiments.

Volumetric Known-Data Attacks: We further evaluate $TWINSSE_{OXT}$ against the known-data volume analysis attacks presented by Blackstone et al. [2], where we analyse $TWINSSE_{OXT}$ against the SelVolAn attack. These specific class of attacks exploits total volume pattern of the queries to recover the keywords, assuming that a fraction of the total data (quantified by “known data rate” δ) is available to the adversary. More precisely, it tries to associate the queried tags available to the server with known keywords. Since $TWINSSE_{OXT}$ produces noisy volumes due to the presence of spurious ids, the recovery rate is expected to be low in these evaluations. We use the LEAKER¹² framework to execute the SelVolAn attack.

We plot the attack results in Figure 13 which depicts the query recovery accuracy for the SelVolAn attack with varying known data rate δ . Note that, as stated earlier, the recovery rate for the SelVolAn attack is significantly low for our construction due the presence of spurious ids resulting in noisy volume pattern. Furthermore, the server receives the volume pattern information of the meta-keywords - not the actual query keywords. From the adversary’s perspective, this requires additional auxiliary information to map recorded encrypted meta-keyword tags to probable meta-keywords pre-computed by the adversary from a partial set of the keyword universe (available as auxiliary data). Precise mapping would require full set of keywords as auxiliary information (indicating a high δ value, an extremely strong assumption) to form actual meta-keywords on the adversary’s side.

Search and Access Pattern (SAP) based Attack: We also evaluate $TWINSSE_{OXT}$ against the state-of-the-art SAP attack by Oya et al. [25]. The SAP attack exploits search pattern (the sequence in which the queries are searched) and the access pattern (the particular address/elements that are “touched” by the server for each queried tag) and combines these two to recover the association among queried tags (recorded by the adversarial server) and the probable keywords available from auxiliary information. Again, in this case too, the server receives noisy access pattern which prevents the adversary to precisely map a recorded tag to a probable keyword (meta-keyword in $TWINSSE_{OXT}$) available as auxiliary data. We validate this through experimental evaluations with $TWINSSE_{OXT}$. We used the code available from the authors¹³ for this evaluation.

The attack results are presented in Figure 14 which depicts the attack accuracy (as a fraction of correct “recorded tag”-“probable

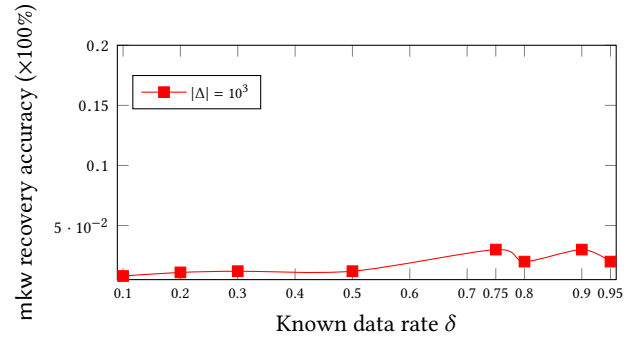


Figure 13: Leakage Analysis of $TWINSSE_{OXT}$ - SelVolAn attack. The amount of information available to the adversary (as a fraction of the total data) is varied and plotted on the x-axis. The volume pattern from meta-keyword queries were supplied as the leaked information.

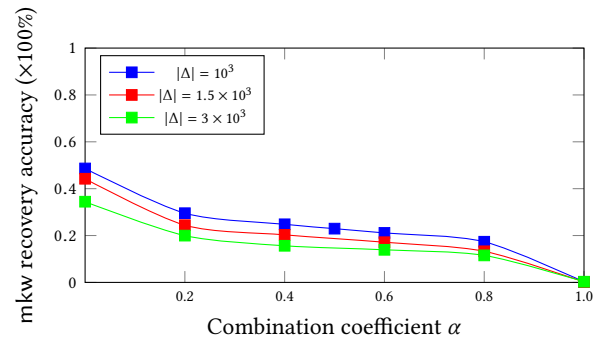


Figure 14: Leakage Analysis of $TWINSSE_{OXT}$ - SAP attack. Access pattern from meta-keyword queries and auxiliary information of keywords from Google Trends were provided as the input data.

meta-keyword” associations) to all such reconstructed associations) with varying combination coefficient (α). At a high-level, α represents the fraction of frequency information of the total auxiliary information available to the adversary used in the attack. In this case, the adversary recovers the association among queried tags and probable meta-keywords - not tags and actual keywords. Since reconstructing the actual meta-keywords requires the exact same keyword universe available to the adversary, it is unlikely that the adversary would be successful in associating the recovered meta-keywords with the subset of the keyword universe available to her (as auxiliary information).

G GENERAL BOOLEAN QUERIES (CNF AND DNF) IN $TWINSSE$ AND $TWINSSE_{OXT}$

In Section 4, we described how $TWINSSE$ and its instantiation from OXT, namely $TWINSSE_{OXT}$, handle purely conjunctive and purely disjunctive queries. In this section, we describe how $TWINSSE$ can be extended to address general Boolean queries in either the conjunctive normal form (CNF) or the disjunctive normal form (DNF).

¹²<https://github.com/encryptogroup/LEAKER>

¹³<https://github.com/simon-oya/USENIX21-sap-code>

We note here that OXT does support Boolean queries beyond simple conjunctions, albeit where the query must be in a restricted *searchable normal form* (SNF) [7]; our transformation is significantly more general in the sense that it extends to any CNF or DNF formula over keywords, well beyond the scope of SNF queries.

We begin by describing how to handle DNF queries, because, similar to its purely conjunctive and disjunctive counterparts, DNF queries are also handled by TWINSSE (and hence, by extension, TWINSSE_{OXT}) by making fully black-box usage of the underlying conjunctive SSE scheme. Subsequently, we show how to address CNF queries. This is slightly more involved, and makes non black-box usage of the underlying conjunctive SSE scheme (we describe a specific strategy for TWINSSE_{OXT} to handle CNF queries that relies on a special data structure used by the OXT scheme).

G.1 Handling Boolean Queries in DNF Form

In Boolean logic, a disjunctive normal form (DNF) is a canonical normal form of a logical formula consisting of a disjunction of conjunctions (alternatively, OR of AND clauses). Formally, any query q that is a Boolean formula over keywords in DNF takes the form

$$q = \bigvee_{\ell \in [L]} q_\ell = \bigvee_{\ell \in [L]} (w_{\ell,1} \wedge \dots \wedge w_{\ell,t_\ell}),$$

where each $q_\ell = (w_{\ell,1} \wedge \dots \wedge w_{\ell,t_\ell})$ for $\ell \in [L]$ is a conjunctive clause. Our approach to handle a DNF query is straightforward, and closely resembles, at a high level, our strategy for handling disjunctive queries via query partitioning in TWINSSE. Let $\text{CSSE} = (\text{CSSE.SETUP}, \text{CSSE.SEARCH})$ be any generic conjunctive SSE scheme. The search algorithm processes q via the following steps (the setup algorithm remains the same as TWINSSE.SETUP described in Algorithm 1, Section 4):

- **Client:** Parse a DNF query as $q = \bigvee_{\ell \in [L]} q_\ell$.
- **Client + Server:** For each $\ell \in [L]$ (either in parallel or in uniformly random order), compute $\text{DB}(q_\ell) = \text{CSSE.SEARCH}(q_\ell, \text{EDB})$, where EDB is the encrypted meta-database output by TWINSSE.SETUP.
- **Client:** Locally compute at the client

$$\text{DB}(q) = \bigcup_{\ell \in [L]} \text{DB}(q_\ell).$$

Correctness. Correctness of search follows immediately from the correctness guarantees of the underlying conjunctive SSE scheme CSSE.

Search Complexity. We present an (asymptotic) analysis of the complexity of handling DNF search queries (more concretely, the computational and communication requirements during DNF query processing) when we instantiate TWINSSE using the OXT protocol from [7], i.e., in TWINSSE_{OXT}. Let q be a DNF query of the form

$$q = \bigvee_{\ell \in [L]} q_\ell = \bigvee_{\ell \in [L]} (w_{\ell,1} \wedge \dots \wedge w_{\ell,t_\ell}),$$

where we assume, without loss of generality, that for each $\ell \in [L]$, $w_{\ell,1}$ is the least frequent conjunct in the conjunctive clause q_ℓ . When processing q using TWINSSE_{OXT}, the computational costs (at both the client and the server) as well as the communication

requirements between the client and the server scale linearly as $O(\gamma_{\text{DNF}})$, where

$$\gamma_{\text{DNF}} = \sum_{\ell \in [L]} t_\ell |\text{DB}(\text{mk}w_{\ell,1})|.$$

Note that this is very similar in flavor to the analysis of disjunctive search query overheads for TWINSSE_{OXT} in Section 4.

Leakage Analysis. We state the following theorems for the leakage from TWINSSE and TWINSSE_{OXT} when processing Boolean queries in DNF form.

THEOREM G.1 (DNF QUERY PROCESSING IN TWINSSE). *Assuming that CSSE is an (adaptively) secure SSE scheme with respect to the leakage function $\mathcal{L}_{\text{CSSE}} = (\mathcal{L}_{\text{CSSE}}^{\text{SETUP}}, \mathcal{L}_{\text{CSSE}}^{\text{SEARCH}})$, the leakage incurred by TWINSSE when processing a DNF query as described above is $\mathcal{L}_{\text{TWINSSE}}^{\text{SEARCH, DNF}}(q)$, where for any DNF query $q = \bigvee_{\ell \in [L]} q_\ell$, we have*

$$\mathcal{L}_{\text{TWINSSE}}^{\text{SEARCH, DNF}}(q) = \left\{ \mathcal{L}_{\text{CSSE}}^{\text{SEARCH}}(q_\ell) \right\}_{\ell \in [L]}.$$

THEOREM G.2 (DNF QUERY PROCESSING IN TWINSSE_{OXT}). *The leakage incurred by TWINSSE_{OXT} when processing a DNF query as described above is $\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH, DNF}}(Q)$, where for any sequence of DNF queries $Q = (q_1, \dots, q_M)$ such that $q_m = \bigvee_{\ell \in [L_m]} q_{m,\ell}$ for each $m \in [M]$, we have*

$$\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH, DNF}}(Q) = [\text{RP}, \text{SP}, \text{EP}, \text{IP}]({q_{m,\ell}}_{\ell \in [L_m]})_{m \in [M]}.$$

where RP, SP, EP and IP leakages for conjunctive queries are as defined in Appendix D.

The proofs of these theorems are very similar to the proofs of Theorems D.1 and D.2 described earlier in Appendix E, and are hence not detailed separately.

G.2 Handling Boolean Queries in CNF Form

In Boolean logic, a conjunctive normal form (CNF) is a canonical normal form of a logical formula consisting of a conjunction of disjunctions (alternatively, AND of OR clauses). Formally, any query q that is a Boolean formula over keywords in the CNF form takes the form

$$q = \bigwedge_{\ell \in [L]} q_\ell = \bigwedge_{\ell \in [L]} (w_{\ell,1} \vee \dots \vee w_{\ell,t_\ell}),$$

where each $q_\ell = (w_{\ell,1} \vee \dots \vee w_{\ell,t_\ell})$ for $\ell \in [L]$ is a disjunctive clause. Our approach to handle a CNF query is slightly more involved, and makes usage of some specific features of the OXT protocol to ensure sub-linear search overheads in practice. Hence, the subsequent description of how to handle CNF queries is specific to TWINSSE_{OXT}. We leave it as an interesting open question to investigate a generic solution using any conjunctive SSE scheme in a black-box manner.

We now describe our proposed strategy for handling CNF queries in TWINSSE_{OXT}. Before delving into the details, we need to recall some details of the original OXT scheme due to Cash *et al.* [7]. We refer the reader to [7] for details of the OXT scheme; however, we will try to make the description here as self-contained as possible. The OXT protocol maintains on the server (as part of the encrypted database EDB) a special data structure called a “cross-tag set” (or XSet in short). The XSet consists of several “cross-tags”, where each

cross-tag $x\text{tag}_{id,w}$ corresponds to a document identity-keyword pair (id, w) , where

$$x\text{tag}_{id,w} \in X\text{Set} \text{ if and only if } w \in \mathbf{DB}(id).$$

In our handling of CNF queries in $\text{TWINSSE}_{\text{OXT}}$, we make black-box usage the following sub-functions provided by any implementation of the OXT protocol:

- $\text{OXT.GENXTAG}(sk, id, w)$: The client can use the secret key generated at setup by OXT.SEARCH to generate $x\text{tag}_{id,w}$ for any document identifier id and keyword w .
- $\text{OXT.SEARCHXTAG}(x\text{tag}_{id,w}; X\text{Set})$: On receipt of a cross-tag $x\text{tag}_{id,w}$ from the client, the server can look up the XSet efficiently to return a bit $\beta \in \{0, 1\}$, where $\beta = 1$ if $x\text{tag}_{id,w} \in X\text{Set}$, and $\beta = 0$ otherwise.

Given these sub-routines, our proposal for processing a CNF query q proceeds via the steps outlined below (the setup algorithm again remains the same as TWINSSE.SETUP described in Algorithm 1, Section 4, albeit for $\text{CSSE} = \text{OXT}$). Note that unlike purely conjunctive/disjunctive queries and DNF queries, all of which required a single round search protocol, our processing of CNF queries now requires two rounds of communication between the client and the server.

- **Client:** Parse a CNF query as

$$q = \bigwedge_{\ell \in [L]} q_\ell = \bigwedge_{\ell \in [L]} (w_{\ell,1} \vee \dots \vee w_{\ell,t_\ell}),$$

- **Client:** Identify the candidate disjunctive clause q_ℓ with the smallest result set (this can be computed in a straightforward manner from the client state st output by OXT.SETUP , which has the frequency of each keyword in the dictionary).
- **Client+Server (Round-1):** Compute the result-set corresponding to the disjunctive clause q_ℓ as

$$\mathbf{DB}(q_\ell) = \text{TWINSSE}_{\text{OXT}}.\text{SEARCH}(q_\ell, \widehat{\mathbf{EDB}}),$$

where $\widehat{\mathbf{EDB}}$ is the encrypted meta-database output by $\text{TWINSSE}_{\text{OXT}}.\text{SETUP}$, by directly using the disjunctive search protocol described in Algorithm 3 (Section 4) with $\text{CSSE} = \text{OXT}$.

- **Client:** For each $id \in \mathbf{DB}(q_\ell)$ and each $w_{i,\ell'}$ for $\ell' \neq \ell$ in the query q , compute

$$x\text{tag}_{id,w_{i,\ell'}} = \text{OXT.GENXTAG}(sk, id, w_{i,\ell'}).$$

- **Client+Server (Round-2):** For each $id \in \mathbf{DB}(q_\ell)$ and each $w_{i,\ell'}$ (either in parallel or in uniformly random order), the client sends $x\text{tag}_{id,w_{i,\ell'}}$ to the server and receives in response

$$\beta_{id,w_{i,\ell'}} = \text{OXT.SEARCHXTAG}(x\text{tag}_{id,w_{i,\ell'}}; X\text{Set}).$$

- **Client:** For each $id \in \mathbf{DB}(q_\ell)$, compute

$$\beta_{id} = \bigwedge_{\ell' \in [L] \setminus \{\ell\}} (\beta_{id,w_{\ell,1}} \vee \dots \vee \beta_{id,w_{\ell,t_{\ell'}}}).$$

Output the final result set

$$\mathcal{R}_q = \{id \in \mathbf{DB}(q_\ell) \text{ such that } \beta_{id} = 1\}.$$

Correctness. Correctness of search follows immediately from the correctness guarantees of $\text{TWINSSE}_{\text{OXT}}$ (Theorem 4.1), and the correctness guarantees of the OXT protocol itself.

Search Complexity. We now present an (asymptotic) analysis of the complexity of handling CNF search queries (more concretely, the computational and communication requirements during CNF query processing) in $\text{TWINSSE}_{\text{OXT}}$. Let q be a CNF query of the form

$$q = \bigwedge_{\ell \in [L]} q_\ell = \bigwedge_{\ell \in [L]} (w_{\ell,1} \vee \dots \vee w_{\ell,t_\ell}),$$

where we assume, without loss of generality, that $q_1 = (w_{1,1} \vee \dots \vee w_{\ell,2})$ is the disjunctive clause with the smallest result set. Let $q_{mkw} = \bigvee_{k \in [n_B]} q_{mkw,k}$ be the corresponding meta-query when the disjunctive search query corresponding to q_1 is processed using $\text{TWINSSE}_{\text{OXT}}.\text{SEARCH}$, and assume without loss of generality that $mkw_{i_k,j_k}^{(k)}$ is the least frequent meta-keyword within $q_{mkw,k}$ for each $k \in [n_B]$ (such that $q_{mkw,k}$ is non-empty).

When processing q using $\text{TWINSSE}_{\text{OXT}}$, the computational costs (at both the client and the server) as well as the communication requirements between the client and the server scale linearly as $O(\gamma_0 + \gamma_1)$,

$$\gamma_0 = \sum_{k \in [n_B]} |q_k| |\mathbf{DB}(mkw_{i_k,j_k}^{(k)})|,$$

where $|q_k|$ denotes the number of meta-keywords in the conjunctive sub-meta-query q_k ($|q_k| = 0$ when q_k is empty), and

$$\gamma_1 = |\mathbf{DB}(q_1)| \cdot (\ell \in [2, L] t_\ell).$$

Note that the term γ_0 is computed exactly as in the analysis of disjunctive search query overheads for $\text{TWINSSE}_{\text{OXT}}$ in Section 4. Moreover, the term γ_1 , which represents computational and communication complexities incurred as a result of the round-2 of the CNF query processing (using OXT.GENXTAG and OXT.SEARCHXTAG), is independent of the frequencies of any of the disjunctive clauses other than the “least frequent clause” q_1 .

Leakage Analysis. We state the following theorems for the leakage from $\text{TWINSSE}_{\text{OXT}}$ when processing Boolean queries in CNF form.

THEOREM G.3 (DNF QUERY PROCESSING IN $\text{TWINSSE}_{\text{OXT}}$). *The leakage incurred by $\text{TWINSSE}_{\text{OXT}}$ when processing a CNF query as described above is $\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH,CNF}}$, where for any sequence of CNF queries $\mathcal{Q} = (q_1, \dots, q_M)$ such that*

$$q_m = \bigwedge_{\ell \in [L_m]} q_{m,\ell} = \bigwedge_{\ell \in [L]} (w_{m,\ell,1} \vee \dots \vee w_{m,\ell,t_{m,\ell}})$$

for each $m \in [M]$, with $q_{m,1}$ being (without loss of generality) the least frequent disjunctive clause for each $m \in [M]$, and for any pair of bucketization parameters (n', n_B) , we have

$$\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH, DNF}}(\mathcal{Q}) = \left([\text{RP}, \text{SP}, \text{EP}, \text{IP}](\mathcal{Q}_{mkw}), \mathcal{L}_{x\text{tag}}^* \right).$$

where RP , SP , EP and IP leakages for conjunctive queries are as defined in Appendix D, and where $\mathcal{Q}_{mkw,1}$ is a sequence of (sub-)meta-queries of the form

$$\mathcal{Q}_{mkw} = \{q_{mkw,k,1}\}_{k \in [n_B]},$$

where for each $\ell \in [M]$, we have

$$q_{mkw,\ell} = \left(\bigvee_{k \in [n_B]} q_{mkw,k,1} \right) = \text{GENMQUERY}(q_{\ell,1}, n', n_B),$$

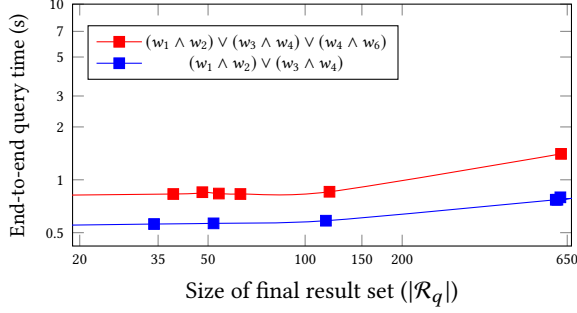


Figure 15: TWINSSE_{OXT} performance with result set size on Enron dataset for DNF queries.

and, finally, we have

$$\mathcal{L}_{\text{xtag}}^* = \{\{\text{DB}(q_{m,1}) \cap \text{DB}(w_{m,\ell}, \ell')\} \mid m \in [M], \ell \in [L_m], \ell' \in [t_{m,\ell}]\}.$$

The proof of this theorem is again very similar to the proof of Theorem D.2 described earlier in Appendix E, and is hence not detailed separately.

G.3 Experimental Results over the Enron Email Dataset

We provide experimental results for CNF and DNF queries using TWINSSE_{OXT} in this section. We experimented over the Enron dataset on the same platform (discussed in Section 5) with our implementation of TWINSSE_{OXT}.

DNF queries. We considered multiple queries with two clauses and three clauses with each clause having two keywords. The end-to-end query time is plotted in Figure 15, where the blue curve represents the query time for two clause queries and the red curve represents the query time for three clause queries. Observe that the query time for both two and three clause queries increase with more number of ids in the final result set. This increment can be attributed to large result size of the individual conjunctive clauses. Also note that the query time increases for three-clause queries due more conjunctive clauses and follows the same trend of increased query time with the final result size.

CNF queries. For experimenting with CNF queries, we considered two-clause queries with two keywords and three keywords per clause. Since the Enron dataset is relatively sparse in nature, with higher number of clauses in query it often results in small or empty intersection. We plotted the end-to-end query time in Figure 16 for both cases – two keyword clauses and three keyword clauses with the size of the final result set. The blue curve represents the end-to-end query time for the queries with two keywords per clause. Similarly, the red curve represents the end-to-end query time for queries with three-keyword clauses. Observe that, in CNF queries also, the end-to-end query time increases with the final result size, due to the increased size of the initial result set. For the three-keyword clauses, the query time is higher than the two-keyword clauses due to the increased size of the initial result set obtained by disjunctive query.

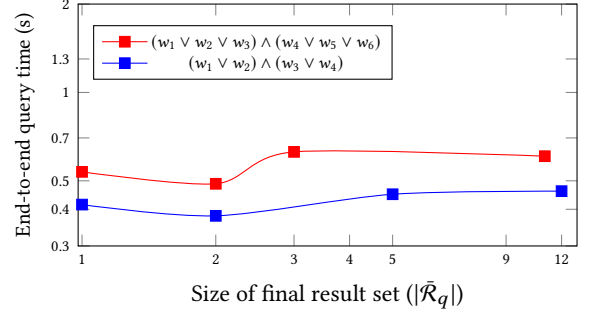


Figure 16: TWINSSE_{OXT} performance with result set size on Enron dataset for CNF queries.

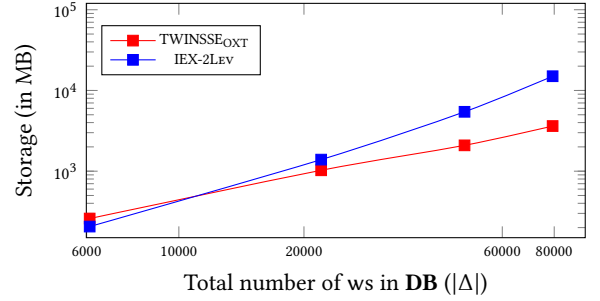


Figure 17: Server storage overhead with plain database size ($|\Delta|$) for the Wikimedia dataset.

H EXPERIMENTAL RESULTS OVER THE WIKIMEDIA DUMP

We present additional experimental results for TWINSSE_{OXT} over Wikimedia databases¹⁴ in this section. We varied the database size from 6K keywords (60K w-id pairs in the plain index) to 80K keywords (8.2 million w-id pairs in the plain index), and we plot the server storage overhead in Figure 17 and performance figures in Figure 18 and Figure 19.

The comparative storage overhead plot (in log scale) in Figure 17 illustrates the quadratic storage overhead for IEX-2LEV; whereas it remains linear for TWINSSE_{OXT}. This storage overhead profile validates our primary contribution of our work, and also illustrates the applicability towards different databases (results on the Enron dataset is presented in the main text Section 5.)

I EVALUATION OF STORAGE OVERHEAD WITH SYNTHETIC DATABASE

We discussed in Section 5 Figure 3 that TWINSSE_{OXT} improves significantly in terms of storage overhead than IEX-2LEV on the Enron database.

Figure 20 compares the storage overhead of TWINSSE_{OXT} and IEX-2LEV on a synthetic database that follows Zipf’s law and Figure 21 compares the estimated storage overhead of TWINSSE_{OXT}

¹⁴<https://dumps.wikimedia.org/enwiki/latest/>

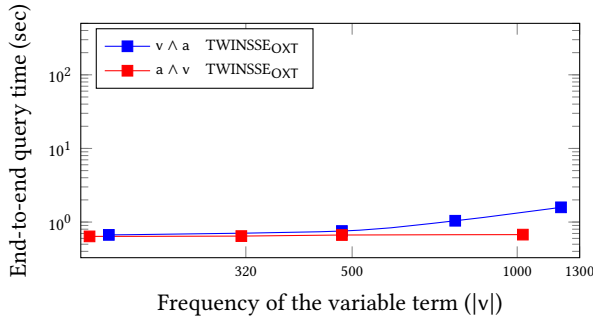


Figure 18: TWINSSE_{OXT} end-to-end conjunctive query search latency vs frequency of the variable term ($|v|$) for Wikimedia dataset.

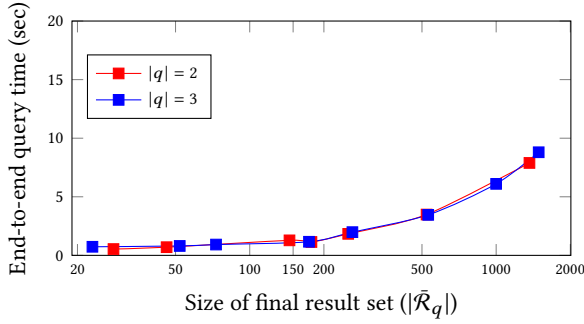


Figure 19: TWINSSE_{OXT} end-to-end disjunctive query search latency vs final result size for Wikimedia dataset.

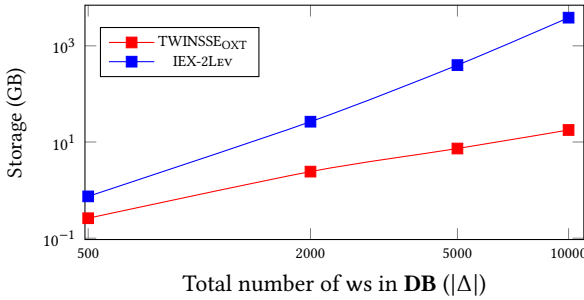


Figure 20: Server storage overhead with number of keywords in synthetic plain database ($|\Delta|$) (prepared following Zipf's distribution).

and IEX-2LEV on a synthetic database that follows a uniform distribution. These databases contain more documents per keyword than the Enron database. This implies that size of the intersections of keyword pairs is much more as compared to the Enron database. Storage overhead of IEX-2LEV hence degrades even more.

To clarify this, the following example of a realistic database can be considered as dense one (as we have described above and in Section 5). Note that, any relational-database is dense if each attribute is low-entropy (i.e., takes only a few values), and hence

each attribute-value pair (equivalent to keywords) occurs in a very large number of records (equivalent to documents). Consider the following Covid-19 patient-database (Table 2), where each attribute-value-pair likely occurs in a large number of patient-records.

Table 2: Example of a dense database with possible attribute-value pairs.

Attributes	Values
Symptomatic	Yes/No
Not-Vaccinated	Yes/No
Dose 1	Yes/No
Dose 2	Yes/No
Booster Dose	Yes/No

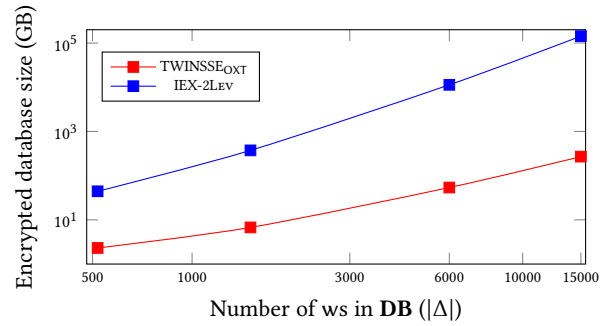


Figure 21: Server storage vs $|\Delta|$ for synthetic DB (following uniform distribution).

Observe that, querying any of the attributes would return a large number of records from this example database. Our experimental results show that IEX-2LEV incurs 70× higher storage overhead than TWINSSE_{OXT} for the synthetic database following Zipf's law (Figure 20) and approximately 150× higher storage overhead for the database following uniform distribution (Figure 21). The search time also increases for both the schemes; however, the main advantage of TWINSSE_{OXT} compared to IEX is in reduced storage, not in search overheads (which still remains sublinear for TWINSSE_{OXT}).