

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

We present the proofs of the lemmas and theorems presented in the main paper here.

## 1 Adaptive Security of SSE

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

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

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

---

```

1: function  $\text{Real}^{\text{SSE}}(\lambda, Q)$ 
2:    $N \leftarrow \text{Adv}(\lambda)$ 
3:    $(\text{sk}, \text{st}_0, \text{EDB}_0) \leftarrow \text{SETUP}(\lambda, N)$ 
4:   for  $k \leftarrow 1$  to  $Q$  do
5:     Let  $q_k \leftarrow \text{Adv}(\lambda, \text{EDB}_{k-1}, \tau_1, \dots, \tau_{k-1})$ 
6:     Let  $(\text{st}_k, \text{EDB}_k, \text{DB}(q_k)) \leftarrow$ 
       SEARCH( $\text{sk}, \text{st}_{k-1}, q_k; \text{EDB}_{k-1}$ )
7:     Let  $\tau_k$  denote the view of the adversary after
       the  $k^{\text{th}}$  query
8:    $b \leftarrow \text{Adv}(\lambda, \text{EDB}_Q, \tau_1, \dots, \tau_Q)$ 
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  $\text{DB}$  other than that encapsulated by  $\mathcal{L}$ . Formally, an SSE scheme is said to be adaptively secure with respect to a leakage function  $\mathcal{L}$  if for any stateful PPT adversary

---

**Algorithm 2** 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  $q_k \leftarrow \mathbf{Adv}(\lambda, \mathbf{EDB}_{k-1}, \tau_1, \dots, \tau_{k-1})$ 
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:    $b \leftarrow \mathbf{Adv}(\lambda, \mathbf{EDB}_Q, \tau_1, \dots, \tau_Q)$ 
9:   return  $b$ 

```

---

$\mathbf{Adv}$  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:

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

where the “real” experiment  $\mathbf{Real}^{\text{SSE}}$  and the “ideal” experiment  $\mathbf{Ideal}^{\text{SSE}}$  are as described in Algorithm 1 and Algorithm 2 (in Appendix).

## 2 Proof of Lemma 3.1 (Already presented in main body)

*Proof.* In the formal proof of Lemma 3.1 we show that, each mkw constructed following the description of Lemma 3.1 covers each  $w$  in  $q$  (the  $\mathbf{DB}(q)$  part) and other  $w$ s that are not in  $q$  are filtered out (the adjacent intersection in Lemma 3.1).

We start with the following conjunctive expression of mkws of  $q_{\text{mkw}}$  for a particular query  $q$  as given in the Lemma 3.1.

$$\begin{aligned}
 \mathbf{DB}(q_{\text{mkw}}) &= \mathbf{DB}(\text{mkw}_{i_1, j_1} \wedge \dots \wedge \text{mkw}_{i_n, j_n}) \\
 &= \bigcap_{k=0}^n \mathbf{DB}(\text{mkw}_{i_k, j_k})
 \end{aligned} \tag{1}$$

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

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

We rewrite Equation (1) in the following way.

$$\begin{aligned}
\mathbf{DB}(q_{\text{mkw}}) &= \bigcap_{k=0}^n \mathbf{DB}(\text{mkw}_{i_k, j_k}) \\
&= \bigcap_{k=0}^n \left( \bigcup_{\ell \in [N] \setminus [l_k, j_k]} \mathbf{DB}(\mathbf{w}_\ell) \right) \\
&= \bigcap_{k=0}^n \left( \bigcup_{r \in [1, n]} \mathbf{DB}(\mathbf{w}_{l_r}) \cup \bigcup_{\substack{\ell \in [n] \setminus (\{[l_k+1, j_k-1]\} \\ \cup \{l_r : r \in [0, n]\})}} \mathbf{DB}(\mathbf{w}_\ell) \right) \\
&= \bigcup_{k \in [0, n]} \mathbf{DB}(\mathbf{w}_{l_k}) \cup \bigcap_{k=0}^n \left( \bigcup_{\substack{\ell \in [N] \setminus (\{[l_k+1, j_k-1]\} \\ \cup \{l_r : r \in [0, n]\})}} \mathbf{DB}(\mathbf{w}_\ell) \right) \\
&= \mathbf{DB}(q) \cup \left( \bigcap_{k=0}^n \left( \bigcup_{\substack{\ell \in [N] \setminus (\{[l_k+1, j_k-1]\} \\ \cup \{l_r : r \in [0, n]\})}} \mathbf{DB}(\mathbf{w}_\ell) \right) \right)
\end{aligned} \tag{2}$$

Observe that, the union inside the right hand term in Equation (2) keeps all  $\mathbf{w}$ s except a stretch of  $\mathbf{w}$ s (from index  $l_k + 1$  to  $l_{k+1} - 1$ ) 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 ids, the following relation holds.

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

We validated this through experimental evaluations (presented in main paper). The intersection on the right of  $\mathbf{DB}(q)$  in the last step of Equation (2) (in Appendix) essentially introduces the spurious *ids* in the result set.  $\square$

### 3 Proof of Theorem 4.1 (Correctness of TWINSSE<sub>BASIC</sub> and TWINSSE)

The proof of correctness for TWINSSE<sub>BASIC</sub> (and TWINSSE as well) follows from the correctness of CSSE. The correctness of CSSE ensures that a conjunctive query  $q = \mathbf{w}_1 \wedge \dots \wedge \mathbf{w}_n$  over an encrypted database satisfies the following relations.

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

$$\mathbf{DB}(\mathbf{w}_1) \cap \dots \cap \mathbf{DB}(\mathbf{w}_n) = \text{CSSE.SEARCH}(q, \mathbf{EDB})$$

We state the proof for TWINSSE<sub>BASIC</sub> first. Then we show that this can be simply extended to main TWINSSE scheme (the final bucketized version).

*Proof for TWINSSE<sub>BASIC</sub>.* Proof of the TWINSSE<sub>BASIC</sub> 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.

$$\mathbf{DB}(q_{\text{mkw}}) = \mathbf{DB}(q) \cup \left( \bigcap_{k=0}^n \left( \bigcup_{\substack{\ell \in [N] \setminus (\{[l_k+1, j_k-1]\} \\ \cup \{l_r : r \in [0, n]\}}}} \mathbf{DB}(w_\ell) \right) \right)$$

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<sub>BASIC</sub>.SEARCH, which proves the correctness of the TWINSSE<sub>BASIC</sub>.

*Proof of TWINSSE.* Recall from Section 4 (in main paper), that in TWINSSE all ws from  $\Delta$  are partitioned into  $n_B$  buckets of uniform size, and we execute the basic meta-keyword generation method developed in TWINSSE<sub>BASIC</sub> over each partition independently. Only those partitions with query meta-keywords are accessed during search.

Assume that the dictionary of ws -  $\Delta$  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  $\Delta_k$  can be expressed in the following way.

$$\Delta_k = \{w_{(k-1)n'+1}, w_{(k-1)n'+2}, \dots, w_{kn'}\}$$

The number of ws in each bucket is denoted by  $n'$ . The set of mkws in each bucket  $\mathcal{W}_k$  are represented by the  $\mathcal{S}_{\text{mkw}, k}$ . The TWINSSE<sub>BASIC</sub> 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 of main paper).

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

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

$$q = \bigvee_{k \in [n_B]} q_k$$

We expand the above expression to individual buckets.

$$\begin{aligned}
\mathbf{DB}(q_{mkw}) &= \bigcup_{k \in [n_B]} \mathbf{DB}(q_{mkw,k}) \\
&= \bigcup_{k \in [n_B]} \left( \mathbf{DB}(q_k) \cup \left( \bigcap_{k=0}^{|q_k|} \left( \bigcup_{\substack{\ell \in [|\mathcal{W}_k|] \setminus \\ (\{[l_k+1, j_k-1]\} \\ \cup \{l_r : r \in [0, |q_k|]\})}} \mathbf{DB}(w_\ell) \right) \right) \right) \\
&= \bigcup_{k \in [n_B]} \mathbf{DB}(q_k) \cup \bigcup_{k \in [n_B]} \left( \bigcap_{k=0}^{|q_k|} \left( \bigcup_{\substack{\ell \in [|\mathcal{W}_k|] \setminus \\ (\{[l_k+1, j_k-1]\} \\ \cup \{l_r : r \in [0, |q_k|]\})}} \mathbf{DB}(w_\ell) \right) \right) \\
&= \mathbf{DB}(q) \cup \bigcup_{k \in [n_B]} \left( \bigcap_{k=0}^{|q_k|} \left( \bigcup_{\substack{\ell \in [|\mathcal{W}_k|] \setminus \\ (\{[l_k+1, j_k-1]\} \\ \cup \{l_r : r \in [0, |q_k|]\})}} \mathbf{DB}(w_\ell) \right) \right)
\end{aligned}$$

Clearly, from the above expression  $\mathbf{DB}(q) \subseteq \mathbf{DB}(q_{mkw})$ , which proves the correctness of results for TWINSSE.

## 4 Detailed Analysis and Discussion on the Leakage of TWINSSE<sub>OXT</sub>

In this section, we first detail the leakage profile for the specific instantiation of TWINSSE based on the OXT scheme, namely TWINSSE<sub>OXT</sub>. We then present a discussion on the leakage profiles for both TWINSSE and TWINSSE<sub>OXT</sub>.

### 4.1 Leakage Profile of TWINSSE<sub>OXT</sub>

In this section, we describe the leakage profile of TWINSSE<sub>OXT</sub>. We begin by recalling from [1] the leakage profile of the original OXT scheme. We then build upon it to describe the leakage profile of TWINSSE<sub>OXT</sub>, 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

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

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  $\text{TWINSSE}_{\text{OXT}}$ .** We now formalize the security of  $\text{TWINSSE}_{\text{OXT}}$  in terms of the leakage profiles described above. We do this using a formal theorem, which may be viewed as a specialization of Theorem 4.2 (in the main paper) to a specific instantiation of  $\text{TWINSSE}$  based on  $\text{OXT}$ . Once again, this theorem is based on the (adaptive) simulation-security definition of SSE in the real world-ideal world paradigm, described in Appendix 1.

**Theorem 4.1** (Security of  $\text{TWINSSE}_{\text{OXT}}$ ).  *$\text{TWINSSE}_{\text{OXT}}$  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  $\mathbf{DB}$ , any sequence of conjunctive queries  $\mathcal{Q}_0 = (q_{1,0}, \dots, q_{M,0})$  and any sequence of disjunctive queries  $\mathcal{Q}_1 = (q_{1,1}, \dots, q_{M',1})$ , and any pair of bucketization parameters  $(n', n_B)$ , we have*

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

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

$$\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH}}(\mathcal{Q}_0, \mathcal{Q}_1) = [\text{RP}, \text{SP}, \text{EP}, \text{IP}](\mathcal{Q}_0, \mathcal{Q}_{\text{mkw},1}),$$

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

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

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 5.

## 4.2 Discussion on the Leakage Profile of $\text{TWINSSE}_{\text{OXT}}$

In this subsection, we present a more in-depth analysis of the leakage profile for  $\text{TWINSSE}_{\text{OXT}}$  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 [2, 1, 3, 4, 5, 6]). This is usually considered acceptable in the SSE literature; indeed the few known data/query recovery attacks that manage to exploit this leakage ([7, 8, 9, 10]) 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 [2, 11, 12, 13, 4, 5, 6] 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, 5]. Fortifying TWINSSE<sub>OXT</sub> 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) [2, 14, 15, 16, 17]. 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<sub>OXT</sub> 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 [9]. More concretely, the adversarial server must be able to compute  $|\mathbf{DB}(w_1) \cap \mathbf{DB}(w_i)|$  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<sub>OXT</sub> 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 [3] that improve upon OXT in terms of provable security against leakage-based cryptanalysis and file-injection attacks.

Finally, fortifying implementations of TWINSSE<sub>OXT</sub> 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<sub>OXT</sub> 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 6, we present a leakage-based cryptanalysis of the TWINSSE<sub>OXT</sub> scheme via experiments over the Enron email corpus. Our experiments help substantiate that the leakages incurred by the disjunctive search protocol in TWINSSE<sub>OXT</sub> 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 [8, 9]. We leave it as an open question to extend the analysis using the more advanced leakage cryptanalysis techniques, such as those proposed in [10, 18].

## 5 Security Proofs for TWINSSE and TWINSSE<sub>EXT</sub>

In this section, we formally prove the security of TWINSSE and TWINSSE<sub>EXT</sub> with respect to the generic and specific leakage profiles described in Theorems 4.2 (in the main paper) and 4.1, respectively.

### 5.1 Proof of Theorem 4.2 (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  $\tau_i$  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  $\text{SIM}_{\text{CSSE}}$  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).

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

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  $q_{\text{mkw}}$ .

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.

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

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.

## 5.2 Proof of Theorem 4.1 (Security Analysis of TWINSSE<sub>OXT</sub>)

We resort to a simulation-based security analysis for TWINSSE<sub>OXT</sub>. 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<sub>OXT</sub>, 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 4.2).

We define the leakage profile of TWINSSE<sub>OXT</sub> 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

$$\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH}}(\widehat{\mathbf{EDB}}, \mathcal{Q}_0, \mathcal{Q}_1) = [RP, SP, EP, IP](\mathcal{Q}_0, \mathcal{Q}_{\text{mkw},1}),$$

For conjunctive queries,

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

Here,  $\{\widehat{RP}, \widehat{SP}, \widehat{EP}, \widehat{IP}\}$  are the  $\{RP, SP, EP, IP\}$  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 TWINSSE<sub>OXT</sub> 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 TWINSSE<sub>OXT</sub>.SEARCH leakage as

$$\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH}}(\widehat{\mathbf{EDB}}, \{q_{mkw,1,k}\}_{k \in [\mathcal{Q}_1]}) = \{\widehat{RP}, \widehat{SP}, \widehat{EP}, \widehat{IP}\}$$

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

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

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

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

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

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

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

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

By the simulation security guarantee of  $OXT$ ,  $TWINSSE_{OXT}$  is 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  $\mathbf{MDB}$  (encrypted to  $\widehat{\mathbf{EDB}}$ ) is abstracted from  $\mathbf{DB}$  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  $\mathbf{DB}$ . 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 5.1, Lemma 5.2, and Lemma 5.3 relate the disjunctive  $q$  with  $w_i \in \Delta$  to the conjunctive  $q$  with  $mkw_i \in \widehat{\Delta}$ .

**Lemma 5.1.** Consider two disjunctive queries of the same length  $t$

$$\begin{aligned} q_0 &= w_{1,q_0} \vee w_{2,q_0} \vee \dots \vee w_{t,q_0}, \quad w_{i,q_0} \\ q_1 &= w_{1,q_1} \vee w_{2,q_1} \vee \dots \vee w_{t,q_1}, \quad w_{i,q_1} \end{aligned}$$

have the following expressions using  $mkws$ ,

$$\begin{aligned} q_0 &= q_{0,mkw} = mkw_{1,q_0} \wedge mkw_{2,q_0} \wedge \dots \wedge mkw_{t+1,q_0} \\ q_1 &= q_{1,mkw} = mkw_{1,q_1} \wedge mkw_{2,q_1} \wedge \dots \wedge mkw_{t+1,q_1} \end{aligned}$$

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 5.1 is given in Section 5.3.1.  $\square$

**Lemma 5.2.** Consider two disjunctive queries  $q_0$  and  $q_1$ , of the same length  $t$  have the mkw expressions as defined in Lemma 5.1 - 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 5.2 is given in Section 5.3.2.  $\square$

**Lemma 5.3.** Consider two disjunctive queries of different length  $t_0$  and  $t_1$  -

$$\begin{aligned} q_0 &= w_{1,q_0} \vee w_{2,q_0} \vee \dots \vee w_{t_0,q_0}, \quad w_{i,q_0} \in \Delta \\ q_1 &= w_{1,q_1} \vee w_{2,q_1} \vee \dots \vee w_{t_1,q_1}, \quad w_{i,q_1} \in \Delta \end{aligned}$$

have following expressions in the mkws

$$\begin{aligned} q_0 &= q_{0,mkw} = mkw_{1,q_0} \wedge mkw_{2,q_0} \wedge \dots \wedge mkw_{t_0+1,q_0} \\ q_1 &= q_{1,mkw} = mkw_{1,q_1} \wedge mkw_{2,q_1} \wedge \dots \wedge mkw_{t_1+1,q_1} \end{aligned}$$

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 5.3 is given in Section 5.3.3.  $\square$

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 5.1, 5.2, and 5.3, an adversary can infer the position of the same ws 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 ws. This is equivalent to the case of two conjunctive queries in keywords having the least-frequent w same.

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

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

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

### 5.3 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.

#### 5.3.1 Proof of Lemma 5.1

*Proof.* By construction, each meta-keyword  $\text{mkw}_i$  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  $\text{mkw}_k$  is same for both the queries  $q_0$  and  $q_1$ . Without loss of generality, a meta-keyword in the basic  $O(N^2)$  ( $\text{TWINSSE}_{\text{BASIC}}$ ) 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  $\text{mkw}$  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  $\text{mkws}$  are constructed using  $w$ s in sorted order, if both queries  $q_0$  and  $q_1$  have the same  $r$  and same  $s$  (as one  $\text{mkw}$  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  $w_{k, q_0} = w_{k, q_1}$ .  $\square$

#### 5.3.2 Proof of Lemma 5.2

*Proof.* We assume the common  $\text{mkw}$  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  $\text{mkw}$  appears at indices  $k_0$  in  $q_0$  and at  $k_1$  in  $q_1$ . Since the indices of  $w$ s in the  $\text{mkw}$  strings are in sorted order (increasing frequency) and remains fixed for all  $\text{mkws}$ , the  $w$ s at index  $r$  and index  $s$  are the same for both  $q_0$  and  $q_1$ . However, as the index of  $\text{mkw}$  is different in  $q_0$  and  $q_1$ , the number of preceding  $w$ s 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  $w_{k_0, q_0} = w_{k_1, q_1}$ .  $\square$

### 5.3.3 Proof of Lemma 5.3

*Proof.* The proof of Lemma 5.3 follows from the proof of Lemma 5.2. Essentially, Lemma 5.3 is the extension of Lemma 5.2 for two different lengths of queries. Intuitively, it can be established in the following way. Recall that in Proof 5.3.2,  $r$  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 5.2 holds.  $\square$

## 6 Leakage Cryptanalysis of TWINSSE<sub>OXT</sub>

In analyzing the leakage profile for TWINSSE<sub>OXT</sub>, we followed a two-step process:

1. **Step 1:** Provably establishing the leakage using a simulation-based framework (see Appendix 5).
2. **Step-2:** Analyzing its impact with respect to well-known attacks such as leakage-abuse and file-injection attacks (please see Appendix 4 for a detailed discussion on why the leakage profile of TWINSSE<sub>OXT</sub> is highly resistant to a large class of existing attacks).

Unfortunately, the SSE literature does not currently have well-defined metrics for *experimentally* evaluating the leakage of a scheme over real-world databases. This is why the two-step approach of leakage enumeration and analysis as mentioned above is accepted widely as the norm in the SSE literature (notably [2, 1, 12, 13, 4, 3, 5]).

In this section, we encapsulate Step-2 of our leakage profile analysis for TWINSSE<sub>OXT</sub> into the following claim:

**Claim 6.1** (Informal). *Conjunctive and disjunctive search operations in TWINSSE<sub>OXT</sub> incur leakages that are almost identical to those incurred by OXT [1] and HXT [3], and are reasonably benign in practice based on analysis already done in prior works.*

We substantiate this claim by leakage cryptanalysis experiments on conjunctive and disjunctive keyword searches in TWINSSE<sub>OXT</sub>. The experiments were conducted over the same Enron email corpus as was used for the performance evaluation experiments in the main paper. Throughout, we use a bucket size  $n' = 10$  (same as for the performance evaluation experiments in Section 5 in main paper).

### 6.1 Leakage Cryptanalysis of Conjunctive Search Queries in TWINSSE<sub>OXT</sub>

We now evaluate the leakage from the conjunctive search protocol in TWINSSE<sub>OXT</sub>. We evaluate the probability that the adversary guesses correctly the keywords  $w_1$

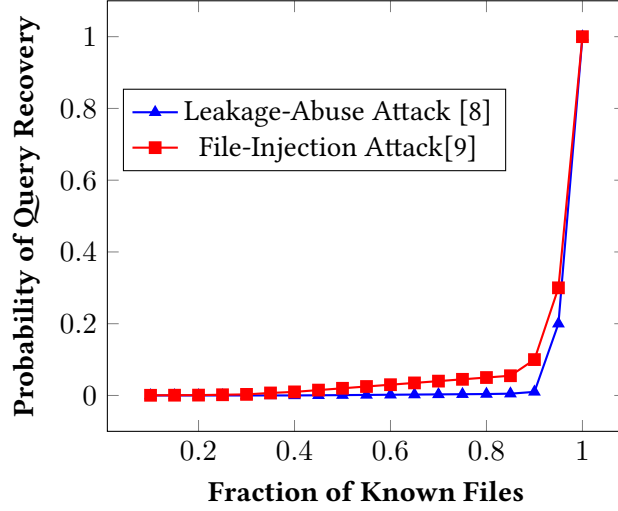


Figure 1: Leakage Analysis of  $\text{TWINSSE}_{\text{OXT}}$ : Two-Keyword Conjunctive Searches in the “Known Files” Setting

and  $w_2$  underlying a two-conjunction query  $q = (w_1 \wedge w_2)$  by one of two well-known and extensively studied cryptanalysis methodologies in the SSE literature—the *leakage-abuse attack* of Cash *et al.* [8] and the *file-injection attack* of Zhang *et al.* [9]. These attacks operate in two models - the *known* file model (where the adversary knows the contents of a certain fraction of the files in the database) and the *chosen/injected* file model (where a certain fraction of the files in the database are adversarially generated).

Naturally, when the adversary knows (or has injected) all the documents in the database, query recovery is trivial. However, this is a very strong attack model and is practically infeasible. What we want in a real-life application is that when the adversary knows only a small fraction of the files in the database, or has managed to inject a small fraction of files into the database, query recovery should happen with a very small probability. This would essentially indicate that the adversary has access to no additional leakage (about either the keywords underlying the query or the files in the database) from the search protocol beyond the benign leakage profile that was formally enumerated in Appendix 5.

Figure 1 illustrates the success probability of the adversary for both kinds of attacks in the “known file” attack setting. The results clearly establish that even when the fraction of known files in the database is as high as 50%, the success probability of the adversary in recovering the keywords underlying a conjunction ( $w_1 \wedge w_2$ ) is less than 5%.

Similarly, Figure 1 illustrates the success probability of the adversary for both kinds of attacks in the “chosen/injected file” attack setting. The results again establish that even when the fraction of injected files in the database is as high as 60% (which is quite unlikely in any real world database), the success probability



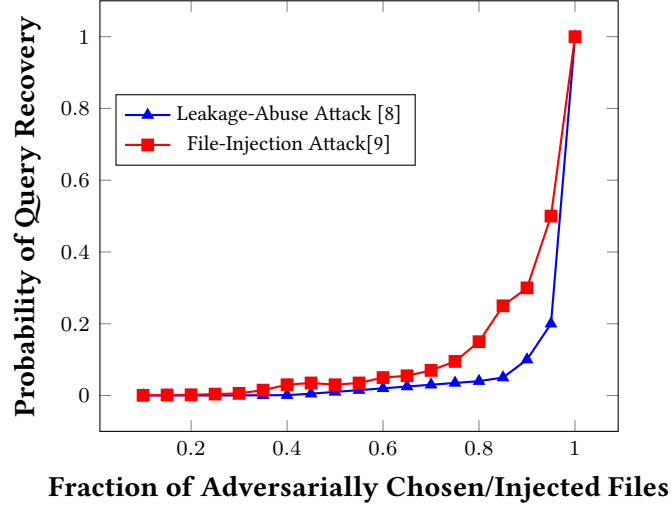


Figure 2: Leakage Analysis of TWINSSE<sub>OXT</sub>: Two-Keyword Conjunctive Searches in the “Chosen Files” Setting

of the adversary in recovering the keywords underlying a conjunction ( $w_1 \wedge w_2$ ) is less than 5%.

Our experiments thus substantiate our statement in Claim 6.1 that the leakages incurred by the disjunctive search protocol in TWINSSE<sub>OXT</sub> 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.

## 6.2 Leakage Cryptanalysis of Disjunctive Search Queries in TWINSSE<sub>OXT</sub>

We now evaluate the leakage from the disjunctive search protocol in TWINSSE<sub>OXT</sub>. We evaluate the probability that the adversary guesses correctly the keywords  $w_1$  and  $w_2$  underlying a two-disjunction query  $q = (w_1 \vee w_2)$  via the same attack methods and the same attack settings in the previous subsection.

Figure 3 illustrates the success probability of the adversary for both kinds of attacks in the “known file” attack setting. The results clearly establish that even when the fraction of known files in the database is as high as 50%, the success probability of the adversary in recovering the keywords underlying a disjunction ( $w_1 \vee w_2$ ) is less than 5%.

Similarly, figure 3 illustrates the success probability of the adversary for both kinds of attacks in the “chosen/injected file” attack setting. The results again establish that even when the fraction of injected files in the database is as high as 60% (which is quite unlikely in any real world database), the success probability of the adversary in recovering the keywords underlying a disjunction ( $w_1 \vee w_2$ ) is less than 5%.

Our experiments additionally substantiate our statement in Claim 6.1 that the leakages incurred by the disjunctive search protocol in TWINSSE<sub>OXT</sub> are reason-

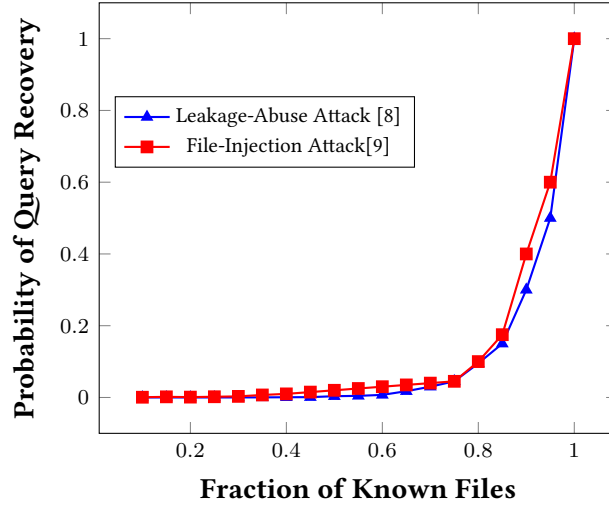


Figure 3: Leakage Analysis of TWINSSE<sub>OXT</sub>: Two-Keyword Disjunctive Searches in the “Known Files” Setting

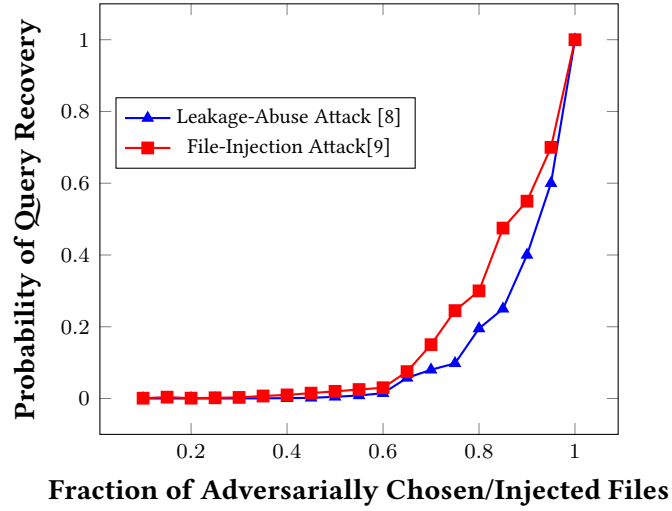


Figure 4: Leakage Analysis of TWINSSE<sub>OXT</sub>: Two-Keyword Disjunctive Searches in the “Chosen Files” Setting

ably benign in practice and are quite resistant to even the most powerful leakage-based cryptanalysis techniques in the SSE literature over real-world databases.

Finally, we leave it as an open question to extend the above analysis using the more advanced leakage cryptanalysis techniques, such as those proposed in [10, 18].

## 7 General Boolean Queries (CNF and DNF) in TWINSSE and TWINSSE<sub>OXT</sub>

In Section 4.3 (in main paper), we described how TWINSSE and its instantiation from OXT, namely TWINSSE<sub>OXT</sub>, 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) form.

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) [1]; 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<sub>OXT</sub>) 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<sub>OXT</sub> to handle CNF queries that relies on a special data structure used by the OXT scheme).

### 7.1 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 the DNF form 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 CSSE = (CSSE.SETUP, 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 (main paper):

- **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  $\widehat{\text{DB}}(q_\ell) = \text{CSSE.SEARCH}(q_\ell, \widehat{\text{EDB}})$ , where  $\widehat{\text{EDB}}$  is the encrypted meta-database output by TWINSSE.SETUP.

- **Client:** Locally compute at the client

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

**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 [1], i.e., in TWINSSE<sub>OXT</sub>. Let  $q$  be a DNF query of the form

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

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  $q_\ell$ . When processing  $q$  using TWINSSE<sub>OXT</sub>, 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

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

Note that this is very similar in flavor to the analysis of disjunctive search query overheads for TWINSSE<sub>OXT</sub> in Section 4 of the main paper.

**Leakage Analysis.** We state the following theorems for the leakage from TWINSSE and TWINSSE<sub>OXT</sub> when processing Boolean queries in DNF form.

**Theorem 7.1** (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  $\mathcal{L}_{\text{TWINSSE}}^{\text{SEARCH, DNF}}$ , where for any DNF query  $q = \bigvee_{\ell \in [L]} q_\ell$ , we have*

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

**Theorem 7.2** (DNF Query Processing in TWINSSE<sub>OXT</sub>). *The leakage incurred by TWINSSE<sub>OXT</sub> when processing a DNF query as described above is  $\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH, DNF}}$ , where for any sequence of DNF queries  $\mathcal{Q} = (q_1, \dots, q_M)$  such that  $q_m = \bigvee_{\ell \in [L_m]} q_{m,\ell}$  for each  $m \in [M]$ , we have*

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

where RP, SP, EP and IP leakages for conjunctive queries are as defined in Appendix 4.

The proofs of these theorems are very similar to the proofs of Theorems 4.2 (in main paper) and 4.1 described earlier in Appendix 5, and are hence not detailed separately.

## 7.2 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  $\text{TWINSSE}_{\text{OXT}}$ . 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  $\text{TWINSSE}_{\text{OXT}}$ . Before delving into the details, we need to recall some details of the original OXT scheme due to Cash *et al.* [1]. We refer the reader to Appendix 10 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  $\text{xtag}_{\text{id},w}$  corresponds to a document identity-keyword pair  $(\text{id}, w)$ , where

$$\text{xtag}_{\text{id},w} \in \text{XSet} \text{ if and only if } w \in \mathbf{DB}(\text{id}).$$

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:

- $\text{OXT.GENXTAG}(\text{sk}, \text{id}, w)$  : The client can use the secret key generated at setup by  $\text{OXT.SEARCH}$  to generate  $\text{xtag}_{\text{id},w}$  for any document identifier  $\text{id}$  and keyword  $w$ .
- $\text{OXT.SEARCHXTAG}(\text{xtag}_{\text{id},w}; \text{XSet})$  : On receipt of a cross-tag  $\text{xtag}_{\text{id},w}$  from the client, the server can look up the XSet efficiently to return a bit  $\beta \in \{0, 1\}$ , where  $\beta = 1$  if  $\text{xtag}_{\text{id},w} \in \text{XSet}$ , 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  $\text{TWINSSE.SETUP}$  described in Algorithm 1 (main paper), 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  $q_\ell$  with the smallest result set (this can be computed in a straightforward manner from the client state st output by  $\text{OXT.SETUP}$ , which has the frequency of each keyword in the dictionary).
- **Client+Server (Round-1):** Compute the result-set corresponding to the disjunctive clause  $q_\ell$  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 (main paper) with  $\text{CSSE} = \text{OXT}$ .

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

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

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

$$\beta_{\text{id},w_{i,\ell'}} = \text{OXT}.\text{SEARCHXTAG}(\text{xtag}_{\text{id},w_{i,\ell'}}; \mathbf{XSet}).$$

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

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

Output the final result set

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

**Correctness.** Correctness of search follows immediately from the correctness guarantees of  $\text{TWINSSE}_{\text{OXT}}$ , and the correctness guarantees of the  $\text{OXT}$  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  $q_1 = (w_{1,2} \vee \dots \vee w_{\ell,2})$  is the disjunctive clause with the smallest result set. Let  $q_{\text{mkw}} = \bigvee_{k \in [n_B]} q_{\text{mkw},k}$  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  $\text{mkw}_{i_k, j_k}^{(k)}$  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  $O(\gamma_0 + \gamma_1)$ ,

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

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), and

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

Note that the term  $\gamma_0$  is computed exactly as in the analysis of disjunctive search query overheads for  $\text{TWINSSE}_{\text{OXT}}$ . Moreover, the term  $\gamma_1$ , which represents computational and communication complexities incurred as a result of the round-2 of the CNF query processing (using  $\text{OXT}.\text{GENXTAG}$  and  $\text{OXT}.\text{SEARCHXTAG}$ ), 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 7.3** (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  $\mathcal{L}_{\text{TWINSSE}_{\text{OXT}}}^{\text{SEARCH}, \text{CNF}}$ , where for any sequence of CNF queries  $\mathcal{Q} = (q_1, \dots, q_M)$  such that*

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

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

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

where  $\text{RP}$ ,  $\text{SP}$ ,  $\text{EP}$  and  $\text{IP}$  leakages for conjunctive queries are as defined in Appendix 4, and where  $\mathcal{Q}_{\text{mkw}, 1}$  is a sequence of (sub-)meta-queries of the form

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

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

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

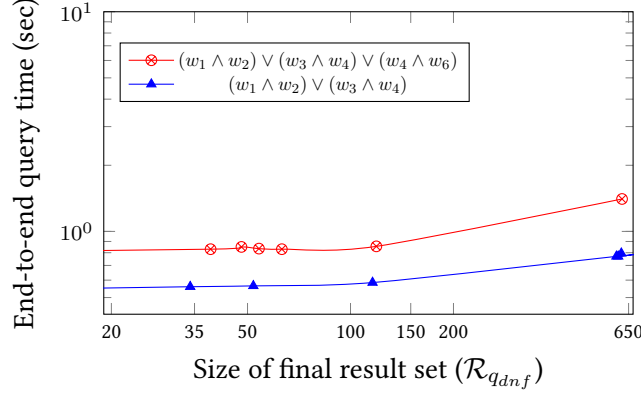


Figure 5: TWINSSE<sub>OXT</sub> performance with result set size on Enron dataset for DNF queries.

and, finally, we have

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

The proof of this theorem is again very similar to the proof of Theorem 4.1 described earlier in Appendix 5, and is hence not detailed separately.

### 7.3 Experimental Results over the Enron Email Dataset

We provide experimental results for CNF and DNF queries using TWINSSE<sub>OXT</sub> in this section. We experimented over the Enron dataset on the same platform (discussed in results of the main paper) with our implementation of TWINSSE<sub>OXT</sub>.

*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 5, 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 6 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



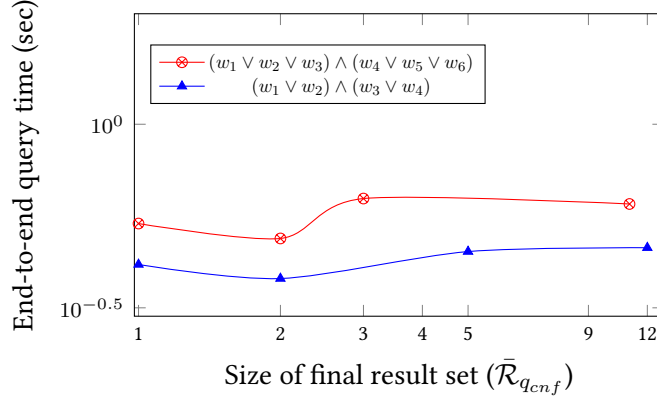


Figure 6: TWINSSE<sub>OXT</sub> performance with result set size on Enron dataset for CNF queries.

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.

## 8 Experimental Results over the Wikimedia Dump

We present additional experimental results for TWINSSE<sub>OXT</sub> over Wikimedia databases<sup>1</sup> 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 7 and performance figures in Figure 8 and Figure 9.

The comparative storage overhead plot (in log scale) in Figure 7 illustrates the quadratic storage overhead for IEX-2LEV; whereas it remains linear for TWINSSE<sub>OXT</sub>. 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 of the main paper.)

## 9 Evaluation of Storage Overhead with Synthetic Database

We discussed in Section 5 Figure 4 (in main paper) that TWINSSE<sub>OXT</sub> improves significantly in terms of storage overhead than IEX-2LEV on the Enron database.

Figure 10 compares the storage overhead of TWINSSE<sub>OXT</sub> and IEX-2LEV on a synthetic database that follows Zipf’s law and Figure 11 compares the estimated storage overhead of TWINSSE<sub>OXT</sub> and IEX-2LEV on a synthetic database

<sup>1</sup><https://dumps.wikimedia.org/enwiki/latest/>

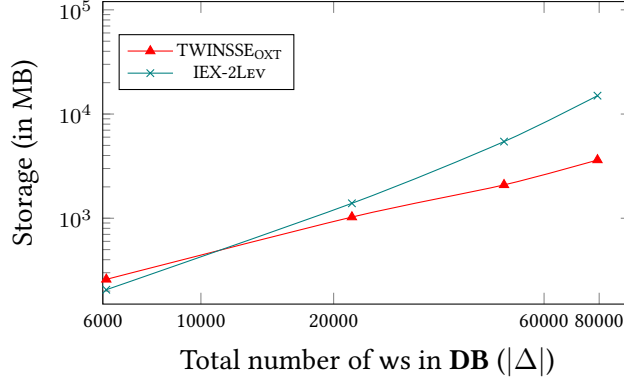


Figure 7: Server storage overhead with plain database size ( $|\mathbf{DB}|$ ) for the Wikimedia dataset.

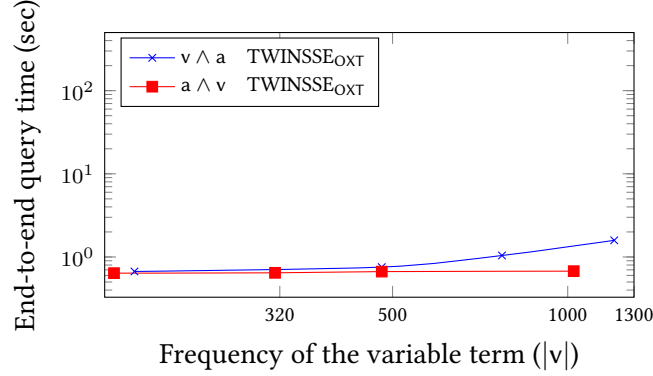


Figure 8: TWINSSE<sub>OXT</sub> end-to-end conjunctive query search latency vs frequency of the variable term ( $|v|$ ) for Wikimedia dataset.

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. Our experimental results show that IEX-2LEV incurs  $70\times$  higher storage overhead than TWINSSE<sub>OXT</sub> for the synthetic database following Zip’s law (Figure 10) and approximately  $150\times$  higher storage overhead for the database following uniform distribution (Figure 11).

## 10 Background Material on OXT

In this section, we present a self-contained description of background material for the OXT scheme from [1].

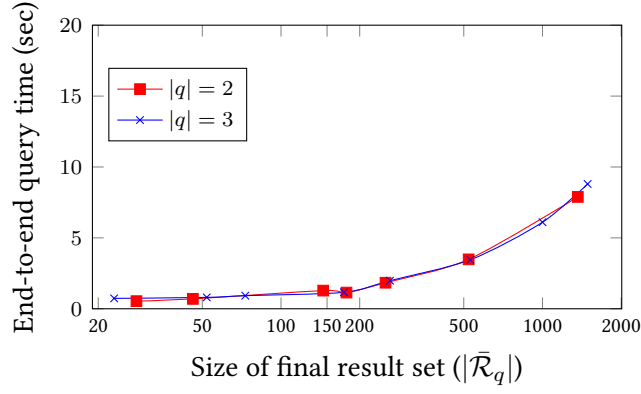


Figure 9: TWINSSE<sub>OXT</sub> end-to-end disjunctive query search latency vs final result size for Wikimedia dataset.

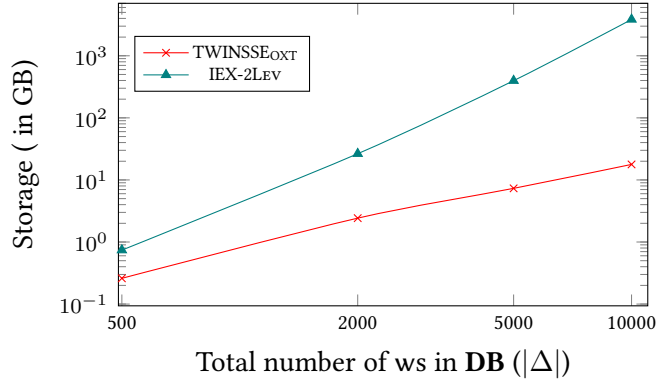


Figure 10: Server storage overhead with number of keywords in synthetic plain database ( $|\Delta|$ ) (prepared following Zipf's distribution).

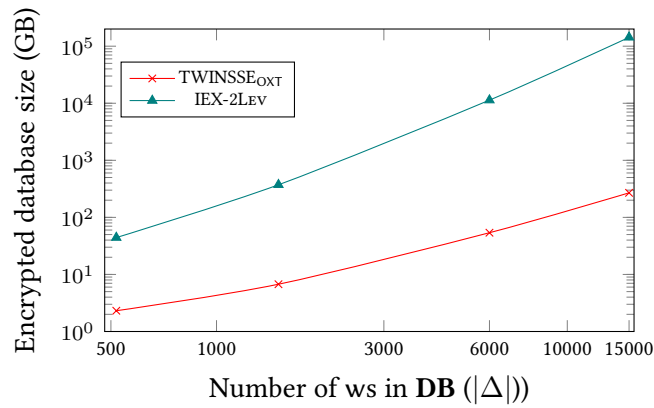


Figure 11: Server storage vs  $|\Delta|$  for synthetic DB (following uniform distribution).

## 10.1 Cryptographic Background

This section presents the definitions and security notions for various cryptographic primitives and assumptions used in the OXT scheme.

**Pseudorandom Functions.** A pseudorandom function (PRF) is a polynomial-time computable function

$$F : \{0, 1\}^\lambda \times \{0, 1\}^\ell \longrightarrow \{0, 1\}^{\ell'},$$

such that for all PPT algorithms  $\mathcal{A}$ , we have

$$\left| \Pr \left[ \mathcal{A}^{F(K, \cdot)} = 1 \right] - \Pr \left[ \mathcal{A}^{G(\cdot)} = 1 \right] \right| \leq \text{negl}(\lambda),$$

where  $K \xleftarrow{R} \{0, 1\}^\lambda$  and  $G$  is uniformly sampled from the set of all functions that map  $\{0, 1\}^\ell$  to  $\{0, 1\}^{\ell'}$ .

**Symmetric-Key Encryption.** A symmetric-key encryption scheme SKE consists of the following polynomial-time algorithms:

- $\text{GEN}(\lambda)$ : A probabilistic algorithm that takes the security parameter  $\lambda$  as input and outputs a secret-key  $\text{sk}$ .
- $\text{ENC}(\text{sk}, x)$ : A probabilistic algorithm that takes as input a key  $\text{sk}$  and a plaintext  $x$ . Outputs a ciphertext  $c$ .
- $\text{DEC}(K, c)$ : A deterministic algorithm that takes as input a key  $\text{sk}$  and a ciphertext  $c$ . Outputs the decrypted plaintext  $x$ .

A symmetric-key encryption scheme is said to be CPA-secure if for all PPT algorithms  $\mathcal{A}$  and any two *arbitrary* plaintext messages  $x_0$  and  $x_1$ , we have

$$\left| \Pr [\mathcal{A}(\text{ENC}(\text{sk}, x_0)) = 1] - \Pr [\mathcal{A}(\text{ENC}(\text{sk}, x_1)) = 1] \right| \leq \text{negl}(\lambda),$$

where  $\text{sk} \leftarrow \text{GEN}(\lambda)$ .

**Decisional Diffie-Hellman Assumption.** Let  $\mathbb{G}$  be a cyclic group of prime order  $p$ , and let  $g$  be any uniformly sampled generator for  $\mathbb{G}$ . The decisional Diffie-Hellman (DDH) assumption is that for all PPT algorithms  $\mathcal{A}$ , we have

$$\left| \Pr \left[ \mathcal{A} \left( g, g^\alpha, g^\beta, g^{\alpha\beta} \right) = 1 \right] - \Pr \left[ \mathcal{A} \left( g, g^\alpha, g^\beta, g^\gamma \right) = 1 \right] \right| \leq \text{negl}(\lambda),$$

where  $\alpha, \beta, \gamma \xleftarrow{R} \mathbb{Z}_p^*$ .

## 10.2 OXT Setup Algorithm

OXT EDB.Setup in Algorithm 3 takes as the security parameter  $\lambda$  and the plaintext database **DB** as input, and generates the encrypted database **EDB** and master key  $\text{mk}$  associated with that database. The encrypted database **EDB** is stored on the remote server, and the client uses the master key to search over the encrypted database. OXT uses a data structure called TSet to store the encrypted data. Detailed construction of TSet is available in OXT paper [1].

---

**Algorithm 3** Setup algorithm of OXT

---

**Input:**  $1^\lambda, \mathbf{DB}$ **Output:**  $\text{mk}, \text{param}, \mathbf{EDB}$ 

```
1: function EDB.SETUP( $1^\lambda, \mathbf{DB}$ )
2:   Initialise  $T \leftarrow \phi$  indexed by keywords  $\Delta$ 
3:   Select key  $K_S$  for PRF  $F$ 
4:   Select keys  $K_I, K_Z, K_X$  for PRF  $F_p$ 
5:   Initialise  $\mathbf{EDB} \leftarrow \{\}$ 
6:   Initialise XSet  $\leftarrow \{\}$ 
7:   for  $w \in \Delta$  do
8:     Initialise  $t \leftarrow \{\}$ 
9:     Compute  $k_e \leftarrow F(K_S, w)$ 
10:    Set counter  $c \leftarrow 1$ 
11:    for  $\text{id} \in \mathbf{DB}(w)$  do
12:      Compute  $x_{\text{id}} \leftarrow F_p(K_I, \text{id})$ 
13:      Compute  $z_w \leftarrow F_p(K_Z, w || c)$ 
14:      Compute  $y_{\text{id}} \leftarrow x_{\text{id}} \cdot z_w^{-1}$ 
15:      Compute  $e_c \leftarrow \text{Sym.Enc}(k_e, \text{id})$ 
16:      Set  $x_{\text{tag}} \leftarrow g^{F_p(K_X, w) \cdot x_{\text{id}}}$ 
17:      Add  $x_{\text{tag}}$  to XSet
18:      Append  $(y_{\text{id}}, e_c)$  to  $t$  and set  $c \leftarrow c + 1$ 
19:    Set  $T[w] \leftarrow t$ 
20:  Compute  $\{\text{TSet}, K_T\} \leftarrow \text{TSet.SetUp}(T)$   $\triangleright$  (See [1] for TSet routines)
21:  return  $\text{mk} = \{\text{msk}, K_S, K_I, K_Z, K_X, K_T\}, \mathbf{EDB} = (\text{TSet}, \text{XSet})$ 
```

---

### 10.3 OXT Search Algorithm

EDB.Search of OXT is a two-party algorithm (Algorithm 4) involving client and hosting server. OXT Search takes the query keywords  $q = \{w_1, \dots, w_n\}$ , master key  $\text{mk}$ , and  $\mathbf{EDB}$  as input and outputs the result of the query  $R_q$ .

---

**Algorithm 4** Search algorithm of OXT

---

**Input:**  $mk, param, q = (w_1 \wedge \dots \wedge w_n), \mathbf{EDB}$

**Output:** Result  $R_q$

```
1: function EDB.SEARCH( $mk, param, q, \mathbf{EDB}$ )
2:   Client's inputs are  $(mk, param, q)$  and server's inputs are  $(param, \mathbf{EDB})$ 
3:   Client initialises  $R_q \leftarrow \{\}$  and computes  $stag \leftarrow \text{TSet.GetTag}(K_T, w_1)$ 
4:   Client sends  $stag$  to the server
5:   Server recovers  $\mathbf{EDB}(1) = \text{TSet}$ . starts accepting  $xtokens$  computed by
      client as follows:
6:   for  $c = 1$  : until server sends stop do
7:     Client computes  $f_{w_1} \leftarrow F_p(k_Z, w_1 || c)$ 
8:     for  $l = 2 : n$  do
9:       Client computes  $xtoken[c, l] \leftarrow g^{f_{w_1} \cdot F_p(k_X, w_l)}$ 
10:      Client sets  $xtoken[c] \leftarrow (xtoken[c, 2], \dots, xtoken[c, n])$ .
11:      Client send  $xtoken[c]$  to server.
12:   Server initialises  $\mathcal{E} \leftarrow \{\}$ .
13:   Server computes  $t \leftarrow \text{TSet.Retrieve}(\text{TSet}, stag)$ ,> (See [1] for TSet rou-
      tines)
14:   for  $c = 1 : |t|$  do
15:     Server recovers  $(y_{id}, e_c)$  from  $c$ -th component of  $t$ .
16:     for  $l = 2 : n$  do
17:       Server computes  $xtag = xtoken[c, l]^{y_{id}}$ 
18:       If  $\forall l \in [2, n], xtag \in \text{XSet}$ , then send  $e_c$  to the client.
19:       When last tuple in  $t$  is reached, send stop to client and halt.
20:   Client computes  $K_e \leftarrow F(K_S, w_1)$ .
21:   Client computes  $id_c \leftarrow \text{Sym.Dec}(K_e, e_c)$ , and adds  $id_c$  to  $R_q$  for all  $e_c$ 
      received.
22:   return  $R_q$ 
```

---