

Efficient Techniques for Privacy-Preserving Sharing of Sensitive Information

Emiliano De Cristofaro, Yanbin Lu, Gene Tsudik
University of California, Irvine

Abstract

The need for controlled (privacy-preserving) sharing of sensitive information occurs in many different and realistic everyday scenarios, ranging from national security to social networking. We consider two interacting parties, at least one of which seeks information from the other: the latter is either willing, or compelled, to share information. This poses two challenges: (1) how to enable this type of sharing such that parties learn no information beyond what they are entitled to, and (2) how to do so efficiently, in real-world practical terms. This paper explores the notion of Privacy-Preserving Sharing of Sensitive Information (PPSSI), and provides two concrete and efficient instantiations, modeled in the context of simple database querying. Proposed techniques function as a *privacy shield* to protect parties from disclosing more than the required minimum of their respective sensitive information. PPSSI deployment prompts several challenges, that are addressed in this paper. Extensive experimental results attest to the practicality of attained privacy features and show that they incur quite low overhead (e.g., 10% slower than standard MySQL).

1 Introduction

In today's increasingly digital world, there is often a tension between safeguarding privacy and sharing information. On the one hand, sensitive data needs to be kept confidential; on the other hand, data owners are often motivated or forced to share sensitive information. Consider the following examples:

- *Aviation Safety*: The Department of Homeland Security (DHS) checks whether any passengers on each flight from/to the United States must be denied boarding or disembarkation, based on several secret lists, including the *Terror Watch List* (TWL). Today, airlines surrender their passenger manifests to the DHS, along with a large amount of sensitive information, including credit card numbers [37]. Besides its obvious privacy implications, this modus operandi poses liability issues with regard to (mostly) innocent passengers' data and concerns about possible data loss.¹ Ideally, the DHS would obtain information pertaining *only* to passengers on one of its watch-lists, without disclosing any information to the airlines.
- *Law Enforcement*: An investigative agency (e.g., the FBI) needs to obtain electronic information about a suspect from other agencies, e.g., the local police, the military, the DMV, the IRS, or the suspect's employer. In many cases, it is dangerous (or simply forbidden) for the FBI to disclose the subjects of its investigation. Whereas, the other party cannot disclose its entire data-set and trust the FBI to only extract desired information. Furthermore, FBI requests might need to be *pre-authorized* by some appropriate authority (e.g., a federal judge). This way, the FBI can only obtain information related to authorized requests.
- *Healthcare*: A health insurance company needs to retrieve information about its client from other entities, such as other insurance carriers or hospitals. The latter cannot provide any information on other patients and the former cannot disclose the identity of the target client.

Other examples of sensitive information sharing include collaborative botnet detection [31], where parties share their logs for the sole purpose of identifying common anomalies.

Motivated by above examples, this paper develops the architecture for **Privacy-Preserving Sharing of Sensitive Information (PPSSI)**, and proposes two efficient and secure instantiations that function as a *privacy shield* to protect parties from disclosing more than the required minimum of sensitive information. We model PPSSI in the context of simple database-querying applications with two parties: a *server*, in possession of a database, and a *client*, performing disjunctive equality queries. In terms of one of the examples above, the airline company (the server) has a database with passenger information, while the DHS (the client) poses queries corresponding to the TWL.

¹See [7] for a litany of recent incidents where large amounts sensitive data were lost or mishandled by government agencies.

Intended Contributions. In this paper, we explore the notion of Privacy-Preserving Sharing of Sensitive Information (PPSSI). Our main building blocks are efficient Private Set Intersection (PSI) techniques. During the design of PPSSI, we address several challenges stemming from adapting PSI to realistic database settings. Our extensive experimental evaluation demonstrates that our techniques incur very low overhead: about 10% slower than standard (non privacy-preserving) MySQL. All source code is publicly available.²

Organization. In next section, we introduce PPSSI syntax, along with its privacy requirements, and review PSI definitions. After reviewing related work in Section 3, in Section 4, we discuss the insecurity of a strawman approach obtained with a naïve adaptation of PSI techniques to PPSSI. Then, Section 5 introduces a secure PPSSI approach using a novel database encryption mechanism. Next, in Section 6, we consider another approach geared for very large databases. Section 7 presents our experimental analysis, and Section 8 concludes the paper by discussing future work.

In Appendix A, we illustrate the details of one PPSSI algorithm, while Appendix B and C present the preliminary version of formal security proofs of solutions proposed in Section 5 and 6. Finally, in Appendix D, we report complete details and performance evaluation of all considered Private Set Intersection constructions.

2 Preliminaries

This section introduces Privacy-Preserving Sharing of Sensitive Information (PPSSI), formalizes its privacy requirements, and overviews Private Set Intersection (PSI) – our main building block.

2.1 PPSSI Syntax & Notation

We model PPSSI in the context of simple database querying. In it, a server maintains a database, DB , containing w records with m attributes $(attr_1, \dots, attr_m)$. We denote $DB = \{(R_j)\}_{j=1}^w$. Each record $R_j = \{val_{j,l}\}_{l=1}^m$, where $val_{j,l}$ is R_j 's value for attribute $attr_l$. A client poses simple disjunctive SQL queries, such as:

$$\text{SELECT } * \text{ FROM } DB \text{ WHERE } (attr_1^* = val_1^* \text{ OR } \dots \text{ OR } attr_v^* = val_v^*) \quad (1)$$

As a result of the query, the client gets all records in DB satisfying *where* clause, and nothing else. Whereas, the server learns nothing about any $\{attr_i^*, val_i^*\}_{1 \leq i \leq v}$. We assume that the database schema (format) is known to the client. Furthermore, without loss of generality, we assume that the client only queries searchable attributes.

In an alternative version supporting *authorized queries*, we require the client to receive query authorizations from a mutually trusted offline *Certification Authority* (CA) prior to interacting with the server. That is, the client outputs matching records only if the client holds pertinent authorizations for $(attr_i^*, val_i^*)$.

Our notation is reflected in Table 1. In addition, we use $Enc_k(\cdot)$ and $Dec_k(\cdot)$ to denote, respectively, symmetric key encryption and decryption (under key k). Public key encryption and decryption, under keys pk and sk , are denoted as $E_{pk}(\cdot)$ and $E_{sk}(\cdot)^{-1}$, respectively. $\sigma = \text{Sign}_{sk}(M)$ denotes a digital signature computed over message M using secret key sk . Operation $\text{Vrfy}_{pk}(\sigma, M)$ returns 1 or 0 indicating whether σ is a valid signature on M . \mathbb{Z}_N^* refers to a composite-order RSA group, where N is the RSA modulus. We use d to denote RSA private key and e to denote corresponding public key. We use \mathbb{Z}_p^* to denote a cyclic group with a subgroup of order q , where p and q are large primes, and $q|p-1$. We use $\hat{e} : G_1 \times G_2 \rightarrow G_t$ to denote a bilinear map and $ZKPK$ to denote zero-knowledge proof of knowledge. We use $H(\cdot), H_1(\cdot), H_2(\cdot), H_3(\cdot)$ to denote different hash functions. In practice, we implement $H(m), H_1(m), H_2(m), H_3(m)$ as $\text{SHA-1}(0||m), \text{SHA-1}(1||m), \text{SHA-1}(2||m), \text{SHA-1}(3||m)$.

2.2 Privacy Requirements

We now define PPSSI privacy requirements for both standard and authorized queries. We consider both Honest-but-Curious (HbC) adversaries and malicious adversaries. An HbC adversary faithfully follows all protocol's specifications (but might attempt to infer additional information during or after protocol execution). Whereas, malicious adversaries may arbitrarily deviate from the protocol.

Privacy requirements are as follows:

- *Server Privacy.* The client learns no information about any record in server's database that does not satisfy the *where* $(attr_i^* = val_i^*)$ clause(s).

²Source code is available at <http://ppssi.googlecode.com/files/ppssi-imp.tar.gz>.

$attr_l$	l th attribute in the database schema	$ctr_{j,l}$	number of times where $val_{j',l} = val_{j,l}, \forall j' \leq j$
R_j	j th record in the database	$tag_{j,l}$	tag for $attr_l, val_{j,l}$
$val_{j,l}$	value in R_j corresponding to $attr_l$	$k'_{j,l}$	key used to encrypt k_j
k_j	key used to encrypt R_j	$k''_{j,l}$	key used to encrypt index j
er_j	encryption of R_j	$ek_{j,l}$	encryption of key k_j
$tk_{j,l}$	token evaluated over $attr_l, val_{j,l}$	$eind_{j,l}$	encryption of index j

Table 1: Notation.

- *Server Privacy (Authorized Queries)*. Same as "Server Privacy" above, but, in addition, the client learns no information about any record satisfying the *where* ($attr_i^* = val_i^*$) clause, unless the ($attr_i^*, val_i^*$) query is authorized by the CA.
- *Client Privacy*. The server learns nothing about any client query parameters, i.e., all $attr_i^*$ and val_i^* , nor about its authorizations, (for authorized queries).
- *Client Unlinkability*. The server cannot determine (with probability non-negligibly exceeding 1/2) whether any two client queries are related.
- *Server Unlinkability*. For any two queries, the client cannot determine whether any record in the server's database has changed, except for the records that are learned (by the client) as a result of both queries.
- *Forward Security (Authorized Queries)*. The client cannot violate Server Privacy with regard to prior interactions, using authorizations obtained later.

Note that Forward Security and Unlinkability requirements are crucial in many practical scenarios. Referring to one example in Section 1, suppose that the FBI queries an employee database without having authorization for a given suspect, e.g., Alice. Server Privacy (Authorized Queries) ensures that the FBI does not obtain any information about Alice. However, unless Forward Security is guaranteed, if the FBI later obtains authorization for Alice, it could inappropriately recover her file from the (recorded) protocol transcript. On the other hand, Unlinkability keeps one party from noticing changes in other party's input. In particular, unless Server Unlinkability is guaranteed, the client can always detect whether the server updates its database between two interactions. Unlinkability also minimizes the risk of privacy leaks. Without Client Unlinkability, if the server learns that the client's queries are the same in two interactions and one of these query contents are leaked, the other query would be immediately exposed.

2.3 Private Set Intersection (PSI)

Private Set Intersection (PSI) [18] constitutes our main building block. It allows two parties – a server and a client – to interact on their respective input sets, such that the client only learns the intersection of the two sets, while the server learns nothing beyond client's set size. Below, we overview two recent PSI protocols [14, 26]. They involve linear communication and computational complexity (in the size of client and server sets).

PSI with Data Transfer (PSI-DT): It involves a server, on input a set of w items, each with associated data record, $\mathcal{S} = \{(s_1, data_1), \dots, (s_w, data_w)\}$, and a client, on input of a set of v items, $\mathcal{C} = \{c_1, \dots, c_v\}$. It results in the client outputting $\{(s_j, data_j) \in \mathcal{S} \mid \exists c_i \in \mathcal{C} \text{ s.t. } c_i = s_j\}$ and the server – nothing except v . This variant is useful whenever the server holds a set of records, rather than a simple set of elements.

Authorized PSI-DT (APSI-DT): It ensures that client input is *authorized* by a mutually trusted offline CA. Unless it holds pertinent authorizations, the client does not learn whether its input is in the intersection. At the same time, the server does not learn whether client's input is authorized, i.e., verification of client authorizations is performed obliviously. More specifically, APSI-DT involves a server, on input of a set of w items: $\mathcal{S} = \{(s_1, data_1), \dots, (s_w, data_w)\}$, and a client, on input of a set of v items with associated authorizations (typically, in the form of digital signatures), $\mathcal{C} = \{(c_1, \sigma_1), \dots, (c_v, \sigma_v)\}$. It results in client outputting $\{(s_j, data_j) \in \mathcal{S} \mid \exists (c_i, \sigma_i) \in \mathcal{C} \text{ s.t. } c_i = s_j \wedge \text{Vrfy}_{pk}(\sigma_i, c_i) = 1\}$ (where pk is CA's public key).

We also distinguish between (A)PSI-DT protocols based on whether or not they support *pre-distribution*:

(A)PSI-DT with pre-distribution: The server can "pre-process" its input set independently from client input. This way, the server can *pre-distribute* its (processed) input before protocol execution. Both pre-processing and pre-distribution can be done offline, once for all possible clients.

(A)PSI-DT without pre-distribution: The server cannot pre-process and pre-distribute its input.

Note that pre-distribution precludes Server Unlinkability, since server input is assumed to be fixed. Similarly, in the context of authorized protocols with pre-distribution, Forward Security cannot be guaranteed.

3 Related Work

A number of cryptographic primitives provide privacy properties resembling those listed in Section 2.2. We overview them below.

Secure Two-Party Computation (2PC). 2PC allows two parties, on input x and y , respectively, to privately compute the output of a public function f over (x, y) . Both parties learn nothing beyond what can be inferred from the output of the computation. Although one could implement PPSSI with generic 2PC, it is usually far more efficient to have dedicated protocols, as 2PC incurs high computational overhead and involves several communication rounds.

Oblivious Transfer (OT). OT [34] involves a sender holding n secret messages and a receiver willing to retrieve the i -th among sender’s messages. It ensures that the sender does not learn which message is retrieved, and the receiver learns no other message. While the OT functionality somehow resembles PPSSI requirements, note that, in PPSSI, receiver’s inputs are queries, whereas, in OT, they are indices.

Private Information Retrieval (PIR). PIR [9] allows a client to retrieve an item from a server database, (1) without revealing which item it is retrieving, and (2) incurring a communication overhead strictly lower than $O(n)$, where n is the database size. Observe that, in PIR, privacy of server’s database is not protected – the client may receive additional bits of information, besides the records requested. Symmetric PIR (SPIR) [20] additionally offers server privacy, thus achieving OT with communication overhead lower than $O(n)$. However, similar to OT, a client of a symmetric PIR needs to input the index of the desired item in server’s database – an unrealistic assumption for PPSSI. An extension to keyword-based retrieval is known as Keyword-PIR (KPIR) [8]. However, KPIR still does not consider server privacy and it involves multiple rounds of PIR executions.

Searchable Encryption (SE). Symmetric Searchable Encryption (SSE) [38] allows a client to store, on an untrusted server, messages encrypted using a symmetric-key cipher under its own secret key. Later, the client can search for specific keywords by giving the server a trapdoor that does not reveal keywords or plaintexts. Boneh et al. [3] later extended SSE to the public-key setting, i.e., anyone can use client’s public key to encrypt and route messages through an untrusted server (e.g., a mail server). The client can then generate search tokens, based on its private key, to let the server identify messages including specific keywords. We conclude that Searchable Encryption targets related yet different scenarios compared to PPSSI.

Privacy-Preserving Database Query (PPDQ). PPDQ techniques can be distinguished into two kinds. The first one is similar to SSE: the client encrypts its data, outsources encrypted data to an untrusted service provider (while not maintaining copies), and queries the service provider at will. In addition to simple equality predicates supported by SSE, solutions like [21, 24] support general SQL operations. Again, this setting is different from PPSSI, as that data, although stored by the server, belongs to the client; thus, there is no privacy restriction against the client. Moreover, these solutions do not provide provably-secure guarantees, but are based on statistical methods.

The second kind of PPDQ is closely related to private predicate matching. Olumofin and Goldberg [33] propose a transition from block-based PIR to SQL-enabled PIR. As opposed to PPSSI, however, server’s database is assumed to be public, thus, its privacy is not protected. Then, Murat and Chris [27] consider a scenario where client matches classification rules against server’s database. However, they assume the client’s rule set to be fixed in advance and known to the server. Additional work, such as [35, 10], requires several independent, mutually-trusted, and non-colluding parties. Murugesan et al. [30] also allow “fuzzy” matching, yet their solution requires a number of (expensive) cryptographic operations (i.e., public-key homomorphic operations) quadratic in the size of parties’ inputs, while we aim at constructing scalable solutions with linear complexity.

4 A Strawman Approach

Looking at definitions in Section 2.3, it seems that PPSSI can be realized by simply instantiating PSI-DT protocols (or APSI-DT for authorized queries). We outline this *strawman* approach below and show that it is not secure.

For each record, consider the hash of every attribute-value pair $(attr_l, val_{j,l})$ as a set element, and R_j as its

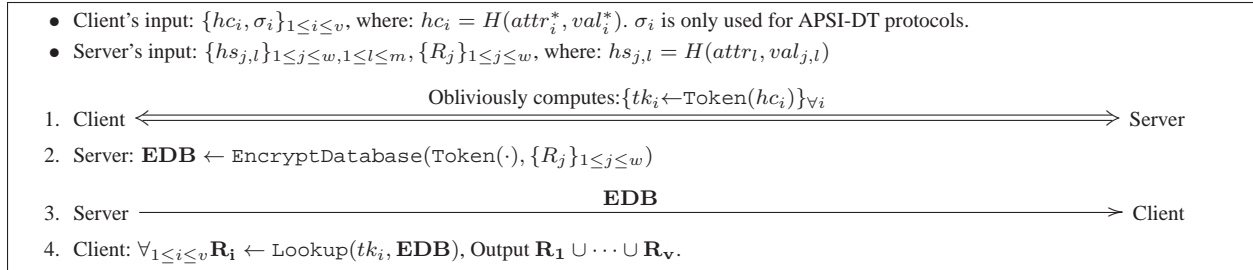


Figure 1: Outline of our first PPSSI approach.

associated data. Server “set” then becomes:

$$\mathcal{S} = \{(H(attr_l, val_{j,l}), R_j)\}_{1 \leq l \leq m, 1 \leq j \leq w}$$

Client “set” is: $\mathcal{C} = \{H(attr_i^*, val_i^*)\}_{1 \leq i \leq v}$, i.e., elements corresponding to the *where* clause in Equation 1. Optionally, if authorized queries are enforced, \mathcal{C} is accompanied by signatures σ_i over $H(attr_i^*, val_i^*)$, following the APSI-DT syntax. Parties engage in an (A)PSI-DT interaction; at the end of it, the client obtains all records matching its query.

The strawman approach faces two security issues:

Challenge 1: Multi-Sets. While most databases include duplicate values (e.g., “gender=male”), PSI-DT and APSI-DT definitions assume that sets do not include duplicates.³ If server set contains duplicated values, the corresponding messages to the client would be identical and the client would learn all patterns and distribution frequencies. This raises a serious concern, as actual values can be often inferred from their frequencies. For example, consider a large database where one attribute reflects “employee blood type”: since blood type frequencies are well-known for general population, distributions for this attribute would essentially reveal the plaintext.

Challenge 2: Data Pointers. To enable querying by any attribute, each record – R_j – must be separately encrypted m times, i.e., once for each attribute. As this would result in high storage/bandwidth overhead, one could encrypt each R_j with a unique symmetric key k_j and then using k_j (instead of R_j) as data associated with $H(attr_l, val_{j,l})$. Although this would reduce the overhead, it would trigger another issue: in order to use the key – rather than the actual record – as the associated “data” in the (A)PSI-DT protocol, we would need to store a pointer to the encrypted record alongside each $H(attr_l, val_{j,l})$. This would allow the client to identify all $H(attr_l, val_{j,l})$ corresponding to a given encrypted record by simply identifying all $H(attr_l, val_{j,l})$ with associated data pointers equal to the given records. Such a (potential) privacy leak would be aggravated if combined with the previous “attack” on multi-sets: given two encrypted records, the client could establish their similarity based on the number of equal attributes.

5 The First PPSSI Approach

We now present our PPSSI construction that is both secure and reasonably practical. Like the strawman approach, it relies on (A)PSI-DT. However, it addresses aforementioned challenges by introducing a novel database-encryption technique. In order to guarantee both *Server Unlinkability* and *Forward Security*, we use (A)PSI-DT *without* pre-distribution.

Our approach is illustrated in Figure 1. In step 1, the client and the server engage in the *oblivious* computation of Token function: at the end of it, the client obtains $tk_i = \text{Token}(hc_i)$, where $hc_i = H(attr_i^*, val_i^*)$. Note that the server learns nothing about hc_i or tk_i . Token function is computed using an (A)PSI-DT protocol, thus, different (A)PSI-DTs instantiate it differently.

In step 2, the server runs `EncryptDatabase` procedure – described in Algorithm 1 and discussed in Section 5.1 – and creates the encrypted database, \mathbf{EDB} that is transferred to the client in step 3. Finally, in step 4, the client runs `Lookup` procedure – illustrated in Algorithm 2 and discussed in Section 5.2 – using tk_i tokens over \mathbf{EDB} ; at the end of it, the client obtains the set of records satisfying its query.

³Note that some PSI constructs (e.g., [29]) support multi-sets, however, their performance is not promising as they incur quadratic computational overhead (in the size of the sets), as opposed to more recent (A)PSI-DT protocols with linear complexity (e.g., [26, 14, 12]). Also, they support neither *data transfer* nor *authorization*.

Algorithm 1: EncryptDatabase Procedure.

input : Function $\text{Token}(\cdot)$ and record set $\{R_j\}_{1 \leq j \leq w}$
output: Encrypted Database **EDB**

1: Shuffle $\{R_j\}_{1 \leq j \leq w}$	11: $k'_{j,l} \leftarrow H_2(tk_{j,l} ctr_{j,l});$
2: $maxlen \leftarrow \max$ length among all R_j	12: $k''_{j,l} \leftarrow H_3(tk_{j,l} ctr_{j,l});$
3: for $1 \leq j \leq w$ do	13: $ek_{j,l} \leftarrow \text{Enc}_{k'_{j,l}}(k_j);$
4: Pad R_j to $maxlen$;	14: $eind_{j,l} \leftarrow \text{Enc}_{k''_{j,l}}(j);$
5: $k_j \xleftarrow{\$} \{0, 1\}^{128}$;	15: LTable $_{j,l} \leftarrow (tag_{j,l}, ek_{j,l}, eind_{j,l});$
6: $er_j \leftarrow \text{Enc}_{k_j}(R_j)$;	16: end for
7: for $1 \leq l \leq m$ do	17: end for
8: $hs_{j,l} \leftarrow H(attr_l, val_{j,l})$;	18: Shuffle LTable with respect to j and l ;
9: $tk_{j,l} \leftarrow \text{Token}(hs_{j,l})$;	19: EDB $\leftarrow \{\text{LTable}, \{er_j\}_{1 \leq j \leq w}\}$;
10: $tag_{j,l} \leftarrow H_1(tk_{j,l} ctr_{j,l})$;	

Our protocol can be used with any (A)PSI-DT, however, we use the variants without pre-distribution, since they provide Server Unlinkability and Forward Security. Following a thorough experimental analysis (Appendix D.5), we select the PSI-DT protocol from [14] (denoted as **DT10-1**) and its APSI-DT counterpart from [14] (denoted as **DT10-APSI**) for authorized queries. These protocols were proven secure against HbC adversaries [14]. However, it was later shown that, with very similar overhead, to achieve security against malicious adversaries [13].

Scheme name	Token definition	PSI category
DT10-1 (Figure 3 of [14])	$\text{Token}(hc) = ((\prod_{i=1}^v hc_i) \cdot g^{R_c}) / hc^{R_s} \bmod p$	PSI-DT without pre-distribution
DT10-APSI (Figure 2 of [14])	$\text{Token}(hc) = ((\prod_{i=1}^v \sigma_i)^2 \cdot g^{R_c}) / hc^{R_s} \bmod N$	APSI-DT without pre-distribution

Table 2: Token definition for (A)PSI-DT without pre-distribution

For the sake of completeness, we define Token function for the selected (A)PSI-DT constructions in Table 2. Note that both Token definitions involve random values R_c and R_s contributed by client and server respectively. Therefore, the server can only evaluate Token over its own inputs after step 1. These random values are selected at the beginning of and kept fixed throughout the PPSSI protocol execution. They are chosen independently, for each invocation, in order to guarantee *Server Unlinkability* and *Forward Security*. We present the complete details of Token 's oblivious computation in Appendix A (see Figure 10 and Figure 11).

Compared to the strawman approach, we modified the “encryption” technique: rather than (directly) using a symmetric-key encryption scheme, the EncryptDatabase procedure is invoked.

5.1 Database Encryption with counters

We illustrate EncryptDatabase procedure in Algorithm 1. It takes in input the definition of the Token function, and server’s record set. It consists of two “phases”: (1) *Record-level* and (2) *Lookup-Table* encryptions.

Record-level encryption is relatively trivial (lines 1–6): first, the server shuffles record locations; then, it pads each R_j up to a fixed maximum record size, picks a random symmetric key k_j , and encrypts R_j as $er_j = \text{Enc}_{k_j}(R_j)$.

Lookup-Table (LTable) encryption (lines 8–15) pertains to attribute name and value pairs. It enables efficient lookup and record decryption. In step 8, the server hashes an attribute-value pair and uses the result as input to Token function in step 9. In step 10, we use the concatenation of Token output and a counter, $ctr_{j,l}$, in order to compute the tag $tag_{j,l}$, later used as a lookup tag during client query. We use $ctr_{j,l}$ to denote the index of duplicate value for the l -th attribute. In other words, $ctr_{j,l}$ is the counter of occurrences of $val_{j',l} = val_{j,l}, \forall j' \leq j$. For example, the third occurrence of value “Smith” for attribute “Last Name” will have the counter equal to 3. The counter guarantees that duplicate $(attr, val)$ pairs correspond to different tags, thus addressing Challenge 1. Next, the server computes $k'_{j,l} = H_2(tk_{j,l} || ctr_{j,l})$ and $k''_{j,l} = H_3(tk_{j,l} || ctr_{j,l})$. Note that $k'_{j,l}$ is used for encrypting symmetric key k_j . Whereas, $k''_{j,l}$ is used for encrypting the index of R_j . In step 13, the server encrypts k_j as $ek_{j,l} = \text{Enc}_{k'_{j,l}}(k_j)$. Then, the server encrypts $eind_{j,l} = \text{Enc}_{k''_{j,l}}(j)$. The encryption of index (data pointer) guarantees that the client cannot link two tags belonging to the same record, thus addressing Challenge 2. In step 15, the server inserts each $tag_{j,l}, ek_{j,l}$ and $eind_{j,l}$ into LTable, which is $\{tag_{j,l}, ek_{j,l}, eind_{j,l}\}_{1 \leq j \leq w, 1 \leq l \leq m}$. Next, the server shuffles LTable (step 18). The resulting encrypted database, **EDB**, is composed of LTable and $\{er_j\}_{j=1}^w$ (step 19).

Algorithm 2: Lookup Procedure.

input : Search token tk and encrypted database $\mathbf{EDB} = \{\mathbf{LTable}, \{er_j\}_{1 \leq j \leq w}\}$
output: Matching record set \mathbf{R}

1: $ctr \leftarrow 1$;	6: $k \leftarrow Dec_{k'}(ek_{j,l})$;
2: while $\exists tag_{j,l} \in \mathbf{LTable}$ s.t. $tag_{j,l} = H_1(tk ctr)$ do	7: $R_j \leftarrow Dec_k(er_{j'})$;
3: $k'' \leftarrow H_3(tk ctr)$;	8: $\mathbf{R} \leftarrow \mathbf{R} \cup R_j$;
4: $j' \leftarrow Dec_{k''}(eind_{j,l})$;	9: $ctr \leftarrow ctr + 1$;
5: $k' \leftarrow H_2(tk ctr)$;	10: end while

5.2 Lookup with counters

We now discuss `Lookup` procedure shown in Algorithm 2. It is used by the client to obtain the query result, i.e., to search \mathbf{EDB} for all records that match client’s search tokens.

In step 1, the client initializes a counter to 1. Next, it searches \mathbf{LTable} for tag $tag_{j,l} = H_1(tk||counter)$. If there is a match, the client attempts to recover the record associated with $tag_{j,l}$. To do so, the client needs to locate the associated record: it computes $k'' = H_3(tk||ctr)$ and recovers $j' = Dec_{k''}(eind_{j,l})$. Note that $er_{j'}$ now corresponds to the associated record. To decrypt $er_{j'}$, the client first recovers the key k used to encrypt $er_{j'}$, by computing $k' = H_2(tk||ctr)$ and obtaining $k = Dec_{k'}(ek_{j,l})$. Finally, the client recovers R_j by decryption, i.e., $R_j = Dec_k(er_{j'})$.

5.3 Correctness

Assume that server’s database includes the attribute “gender” with two occurrences of value “male”. In Algorithm 1, the same tk (step 9) will be generated for the two occurrences of (“gender”, ”male”). However, for the first occurrence, $tag = H_1(tk||1), k' = H_2(tk||1), k'' = H_3(tk||1)$ while, for the second occurrence, $tag = H_1(tk||2), k' = H_2(tk||2), k'' = H_3(tk||2)$.

Suppose that the client searches for records matching “gender = male”, it first derives tk (step 1 of Figure 1). Next, it matches $H_1(tk||1)$ in \mathbf{LTable} , derives keys $k' = H_2(tk||1), k'' = H_3(tk||1)$, and recovers the index in step 4 and the record in step 7 of Algorithm 2. It also looks for $H_1(tk||2)$ and performs the same operations as before, except that $k' = H_2(tk||2), k'' = H_3(tk||2)$. Finally, the client looks for $H_1(tk||3)$: since it finds no match, it terminates.

5.4 Challenges Revisited

We claim that our approach addresses Challenge 1 and 2, discussed in Section 4. While we defer formal proofs Appendix B, the intuition is as follows:

Multi-sets: The use of counters during database encryption makes each $tag_{j,l}$ (resp. $ek_{j,l}, eind_{j,l}$) distinct in \mathbf{LTable} , thus hiding plaintext patterns.

Data Pointers: Storing $eind_{j,l}$ (rather than j) in \mathbf{LTable} , prevents the server from exposing the relationship between an entry $\mathbf{LTable}_{j,l}$ and its associated record R_j .

6 The Second PPSSI Approach for Very Large Databases

The first PPSSI approach in Section 5, combines efficiency with provably-secure guarantees. However, in the context of *very large* databases, it faces two additional issues:

Challenge 3: Bandwidth. If server’s database is very large and/or communication takes place over a slow channel, the bandwidth overhead incurred by the transfer of the encrypted database may become prohibitive.

Challenge 4: Liability. The transfer of the encrypted database to the client also prompts the problem of long-term data safety and associated liability. An encryption scheme considered strong today might gradually weaken in the long term. While we ensure that the client cannot decrypt records outside its query, it is not too far-fetched to imagine that the client might decrypt the entire database in reasonably near future, e.g., 10 or 20 years later. However, data sensitivity might not dissipate over time. For example, suppose that a low-level DoD employee is only allowed to access unclassified data. By gaining access to the encrypted database containing top secret data and patiently

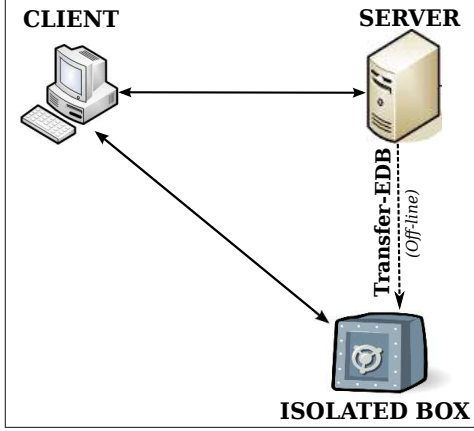


Figure 2: The introduction of the Isolated Box.

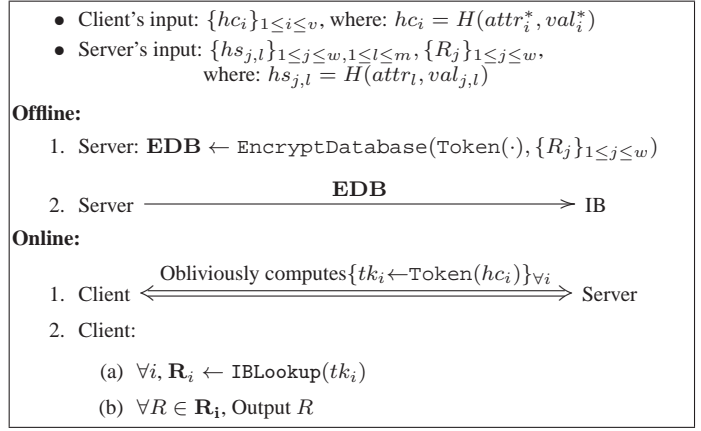


Figure 3: Outline of our second PPSSI approach based on IB.

waiting for the encryption scheme to “age”, the employee might obtain still-classified sensitive information. Further, in several settings, parties (e.g., banks) may be prevented, by regulation, from releasing copies of their databases (even if encrypted).

In the rest of this section, we introduce a novel architecture to address the challenges for very large databases. Our new approach incurs very limited overhead (in terms of both computation and communication), even when compared to non-privacy preserving querying systems.

6.1 Introducing the “Isolated Box”

In order to address Challenge 3 and 4, we propose a system architecture shown in Figure 2. It includes a new component: “Isolated Box” (IB), a non-colluding, untrusted party connected with both the server and the client.

The new interaction involving IB is shown in Figure 3. During the (offline) setup phase, the server encrypts its database, using `EncryptDatabase` (Algorithm 1), and transfers the encrypted database to the IB. Server’s computation of `Token` functionality no longer depends on client’s input, thus, the server can evaluate `Token`(\cdot) without involving the client.

To pose a query, the client first engages with the server in oblivious computation of `Token` (online step 1). Next, for each computed token, it runs the `IBLookup` procedure (Algorithm 3) to retrieve matching records from the IB.

The `Token`(\cdot) functionality is now instantiated using (A)PSI-DT *with* pre-distribution. Specifically, we select the construction from [14] (denoted as **DT10-2**), [26] (denoted as **JL10**) and [12] (denoted as **IBE-APSI**). Again, our choices are based on these protocols’ efficiency and security models. Our experiments – in Appendix D.5 – show that DT10-2, secure in the presence of HbC adversaries, is the most efficient construction, while JL10 combines reasonable efficiency with security against malicious adversary. IBE-APSI is the only APSI-DT with pre-distribution, and it is secure against HbC adversaries. For the sake of completeness, we define `Token` function for the selected (A)PSI-DT constructions in Table 3. Note that d, k, z are server’s secret parameters. Complete details, for each instantiation, are presented in Appendix A.

Scheme name	Token definition	PSI category
DT10-2 (Figure 4 of [14])	$\text{Token}(hc) = (hc)^d \bmod N$	PSI-DT with pre-distribution
JL10 (Figure 2 of [26])	$\text{Token}(hc) = ((hc)^{(p-1)/q})^k \bmod p$	PSI-DT with pre-distribution
IBE-APSI (Figure 5 of [12])	$\text{Token}(hc) = \hat{e}(Q, hc)^z$	APSI-DT with pre-distribution

Table 3: Token for (A)PSI-DT with pre-distribution

Trust Assumptions. The Isolated Box is assumed not to collude with either the server or the client. (Although, we discuss the consequences of collusion in Section 6.6.) We remark that the use of non-colluding parties in the context of Secure Computation was first suggested by [15], and then applied in [28, 27, 1].

Algorithm 3: IBLookup Procedure

Client's input : tk_i
IB's input : $\text{EDB} = \{\text{LTable}, \{er_j\}_{1 \leq j \leq w}\}$
Client's output: Matching record set \mathbf{R}

1. Client: $ctr \leftarrow 1$
2. Client: $tag_i \leftarrow H_1(tk_i || ctr), k_i'' \leftarrow H_3(tk_i || ctr)$
3. Client $\xrightarrow{tag_i, k_i''}$ IB
4. IB: If $(\exists tag_{j,l} \in \text{LTable}_{j,l} \text{ s.t. } tag_{j,l} = tag_i)$
 $j' \leftarrow Dec_{k_i''}(eind_{j,l}),$
 $ret \leftarrow \{ek_{j,l}, er_{j'}\}$
5. IB \xrightarrow{ret} Client
 - else
 $ret \leftarrow \perp$
6. Client: If $ret = \perp$, abort
 else $k_i' = H_2(tk_i || ctr), k_i = Dec_{k_i'}(ek_{j,l})$
 $R_i = Dec_{k_i}(er_{j'}), \mathbf{R} \leftarrow \mathbf{R} \cup R_i$
 $ctr \leftarrow ctr + 1$, Goto step 2.

While our requirement for the presence of IB might seem like a “strong” assumption, we stress that the IB is only trusted not to collude with other parties. It simply stores server’s encrypted database and return ciphertexts matching client’s encrypted queries (i.e., *tags*), without learning any information about records and queries. Also note that, in practice, the IB can be either instantiated as a (non-colluding) cloud server or as a piece of secure hardware installed on server’s premises: it is only important to ensure that the server does not learn *what* the IB reads from its storage and transfers to the client.

6.2 Database Encryption

IB’s presence does not really affect database encryption, i.e., `Encryptdatabase` procedure presented in Algorithm 1. It only uses a different `Token(·)` function. While in the first approach (Section 5) we rely on (A)PSI-DT *without* pre-distribution (i.e., the server cannot run `Token(·)` before interacting with the client), we now use (A)PSI-DT *with* pre-distribution. Thus, the server can evaluate `Token(·)` over its own inputs, *offline*, and then transfer the encrypted database to the IB.

6.3 Query lookup

IBLookup procedure is used by the client to obtain records matching client’s query. It is shown in Algorithm 3. Similar to our first approach, the client runs the lookup procedure after obtaining search tokens (via oblivious computation of `Token` – online step 1 in Figure 3). For each derived token, tk_i , it invokes IBLookup to retrieve (from the IB) all records matching tk_i .

We use the term *transaction* to denote a complete query procedure, for each tk_i (from the time the first query for tk_i is issued, until the last response from the IB is received). *Retrieval* denotes the receipt of a single response record during a transaction. A transaction is composed of several retrievals between the client and the IB. The client retrieves records one by one from the IB, by gradually incrementing the counter ctr . In step 1, the client sets ctr to 1. In step 2, the client derives tag_i and an index decryption key k_i'' from token tk_i . After receiving tag_i and k_i'' in step 3, the IB searches for matching tags in the lookup table in step 4. If there is a match, the IB recovers the index j' by decrypting $eind_{j,l}$ with k_i'' , assembles the corresponding record $er_{j'}$ and the ciphertext of its decryption key $ek_{j,l}$ into ret and transmits ret to the client in step 5. Otherwise, \perp is transmitted. If the client receives \perp , it aborts. Otherwise, it decrypts $ek_{j,l}$ into k_i with k_i' and recovers record R_i from $er_{j'}$ using k_i . Then, it increments ctr and starts another retrieval by returning to step 2.

6.4 Optimizations

Since transmission of ret may incur some delay, Algorithm 3 can be sped up by pipe-lining computation of tag_i and k_i'' (step 2) in next retrieval with the transmission of ret (step 5) in current retrieval.

Note that the computation of $ek_{j,l}$ and $eind_{j,l}$ (steps 13–14 in Algorithm 1) can also be optimized. Since we use a counter as input to compute $k_{j,l}'$ (respectively, $k_{j,l}''$), each $k_{j,l}'$ (respectively, $k_{j,l}''$) is different for any j, l . Both $k_{j,l}'$ and $k_{j,l}''$ are 160-bit values (SHA-1), while k_j is 128 bits and j is clearly smaller. Hence, we can use *one-time-pad* encryption (i.e. $ek_{j,l} = k_{j,l}' \oplus k_j$ and $eind_{j,l} = k_{j,l}'' \oplus j$) to speed up computation. In Algorithm 3, $Dec_{k_i''}(eind_{j,l})$ becomes $k_i'' \oplus eind_{j,l}$ and $Dec_{k_i'}(ek_{j,l})$ changes to $k_i' \oplus ek_{j,l}$.

6.5 Challenges Revisited

Since we use the same encryption procedure discussed in Section 5, Challenge 1 and 2 are already addressed. Thus, we only consider Challenge 3 and 4.

Bandwidth: Once the server transfers its database (offline) to the IB, the latter returns to the client only records matching its query. Therefore, bandwidth consumption is minimized.

Liability: Since the IB holds the encrypted database, the client only obtains the result of its queries, thus, ruling out any potential liability issues.

Finally, the introduction of the IB enables Server Unlinkability and Forward Security, despite the fact that we use (A)PSI-DT *with* pre-distribution techniques. Indeed, records not matching a query are never available to the client, thus, it does not learn whether they have changed. Similarly, the client cannot use future authorizations to maliciously obtain information from previous (recorded) interactions.

6.6 Discussion

Privacy Revisited. The introduction of the IB and the use of counter mode in database encryption provide additional privacy properties. If the client performs only one query transaction, as in Algorithm 3, the IB can link all *tag* values in step 3 to the same $(attr, val)$ pair. This may pose a similar risk to that discussed in the “multi-set” challenge, with respect to the IB. However, the counter allows the client to retrieve matching records one by one. Therefore, the client can choose to add a random delay between two subsequent retrievals in a single transaction. If the distribution of additional delay is indistinguishable from time gaps between two transactions, the IB cannot tell the difference between two continuous retrievals within one transaction from two distinct transactions. As a result, the IB cannot infer whether two continuously retrieved records share the same $(attr, val)$ pair and the distribution of the attribute value remains hidden.

Also note that the introduction of the IB does not violate Client or Server Privacy. Client Privacy is preserved because the client (obviously) computes a token, which is not learned by the server. The IB does not learn client’s interests, since client’s input to the IB (*tag*) is statistically indistinguishable from a random value. Server Privacy is preserved because the client does not gain any extra information by interacting with the IB. Finally, the IB only holds the encrypted database and learns no plaintext.

Removing Online Server. Although it only needs to perform oblivious computation of tokens, we still require the server to be online. Inspired by [22] and [16], we can replace the online server with a tamper-proof smartcard, dedicated to computing Token function. The server only needs to program its secret key into the smartcard, which protects the key from being accessed by the client. This way, after handing the smartcard to the client, the server can go offline. The smartcard is assumed to enforce a limit on the number of Token invocations.

Limitations. We acknowledge that our second PPSSI approach has some limitations. Over time, as it serves many queries, the IB gradually learns the relationship between tags and encrypted records through pointers associated with each tag. This issue can be mitigated by letting the server periodically re-encrypt the database. IB also learns database access patterns generated by query executions. Nonetheless, without knowing the distribution of query predicates, the access pattern of encrypted data leaks very little information to the IB. Next, if the server and the IB collude, Client Privacy is lost, since the IB learns *tag* that the client seeks, and the server knows the $(attr, val)$ pair each *tag* is related to. On the other hand, if the client and the IB collude, the client can access the entire encrypted database, thus, liability becomes a problem. Last, Server Unlinkability is protected only with respect to the client. Server Unlinkability with respect to the IB is not guaranteed, since the IB learns about all changes in server’s database. Finally, note that PPSSI currently supports only equality and disjunctive queries. Enabling conjunctive queries would require treating all combinations of $(attr, val)$ pairs as server’s set elements. Thus, client’s input would become exponential in terms of the number of attributes. This remains an interesting challenge left as part of future work.

7 Performance Evaluation

In this section, we evaluate the performance of our PPSSI approaches. First, we benchmark cryptographic operations and use these results to derive step-by-step cost of proposed techniques. Next, we compare our first PPSSI ap-

proach to PIR. Finally, we build a (limited) DBMS to compare our second PPSSI approach to a non privacy-preserving MySQL database.

7.1 Benchmarking All PPSSI Components

The following benchmark refers to executions on an Intel Harpertown server with Xeon E5420 CPU (2.5 GHz, 12MB L2 Cache) and 8GB RAM inside. We build the benchmarking tool based on OpenSSL library (ver.1.0.0c) and PBC library (ver.0.5.11).

7.1.1 Cryptographic Operations

We start with benchmarking modular arithmetic operations. In Table 4, we present performance results for modular multiplication (**mul**) and modular inversion (**inv**) under different modulus sizes (column 1). We also report the performance of modular exponentiation (**exp**) and modular exponentiation with Chinese Remainder Theorem (**exp.crt**) under different combinations of modulus sizes (column 1) and exponent sizes (column 4). We choose modulus size to be 1024, 2048, 3072 bits respectively, which corresponds to 80, 112, and 128 symmetric key security level. (The protection lifetime of 1024-bit modulus is supposed to last until 2010, whereas, that of 2048-bit modulus is until 2030, and 3072-bit – to 2030 and beyond [5]). **exp.crt** can only be used when factorization of N is known, thus, we only measure its performance for exponent size $|d|$ (being d RSA secret key).

modulus (bits)	mul (ms)	inv (ms)	exponent (bits)	exp (ms)	exp.crt (ms)
$ q =160$	0.001	0.016	–	–	–
$ p =1024$	0.003	0.244	$ q =160$	0.297	–
			$ p =1024$	1.725	–
$ N =1024$	0.003	0.244	$ d =1024$	1.725	0.534
			$ e =17$	0.039	–
$ q =256$	0.001	0.03	–	–	–
$ p =2048$	0.009	0.765	$ q =256$	1.685	–
			$ p =2048$	12.679	–
$ N =2048$	0.009	0.765	$ d =2048$	12.679	3.451
			$ e =17$	0.124	–
$ p =3072$	0.02	0.837	$ q =256$	3.719	–
			$ p =3072$	41.784	–
$ N =3072$	0.02	0.837	$ d =3072$	41.784	11.031
			$ e =17$	0.263	–

Table 4: Benchmarking **mul** and **exp** operations using the OpenSSL library.

Table 5 shows the benchmark results of operations in bilinear map $\hat{e} : G_1 \times G_2 \rightarrow G_t$ under different G_1/G_2 base size and different group orders. We choose type A pairing provided in PBC library. Since type A provides $2 \cdot |base|$ discrete logarithm security, we use half the group size as we do in Table 4. We use **exp**(G_1) and **exp**(G_2) to denote exponentiation in group G_1 and G_2 respectively.

base (bits)	order (bits)	exp in G_1/G_2 (ms)	exp in G_T (ms)	pairing (ms)
512	160	2.492	0.233	1.859
1024	256	8.896	0.998	9.481
1536	256	15.086	1.922	21.826

Table 5: Benchmarking operations on bilinear maps using the PBC library.

In Table 6, we evaluate different symmetric encryption schemes and hash functions. For symmetric encryption, we only experiment with 128-bit key size, since it is the lowest supported by AES and it matches the security level of 3072-bit RSA keys. The decryption cost is same as the encryption cost, hence, we omit it here.

Symmetric encryption (ms/MB)			Hash function (ms/MB)		
RC4	AES-CBC	AES-CTR	SHA1	SHA256	SHA512
3.500	6.539	13.820	3.406	6.867	4.586

Table 6: Benchmarking (128-bit) symmetric-key encryptions and hash function computations.

7.1.2 PPSSI Operations

We now evaluate the performance of all operations involved in both of our PPSSI approaches. Remark that we use 2048-bit modulus and records of fixed $2KB$ length.

Figure 4 measures the time needed to perform the oblivious computation of `Token` function, for every possible (A)PSI-DT instantiation. Observe that the cost always increases linearly with client’s query size. As for protocols without pre-distribution, DT10-APSI is unsurprisingly more expensive than DT10-1. Whereas, DT10-2 and JL10 are, respectively, the most and the least efficient ones of protocols with pre-distribution.

Then, Figure 5 evaluates the performance of the Lookup-Table encryption, performed by the server. This operation includes server’s computation of `Token` function over its own input (Note that this is not oblivious computation). Again, running time always increase linearly with the product of the number of records (w) and the number of attributes (m).

In Figure 6, we report the cost of the Record-level encryption. This only depends on the number of records. Compared to the Lookup-table encryption, the Record-level encryption incurs a negligible overhead.

Finally, Figure 7 presents the running time of the Lookup procedures (Algorithms 2 and Algorithms 3 without consideration of communication delay). Unsurprisingly, cost is identical for both algorithms and increases linearly with the number of matching records (v_m).

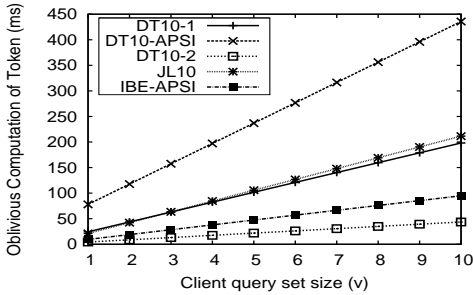


Figure 4: Token Oblivious Computation.

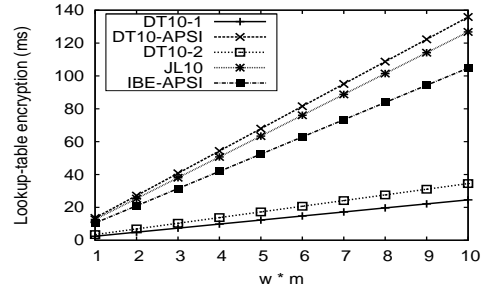


Figure 5: Lookup-Table Encryption (line 8-15 of Algorithm 1).

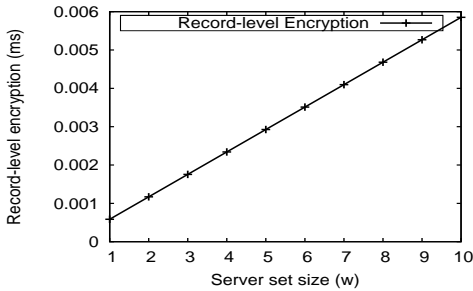


Figure 6: Record-level Encryption (line 1-6 of Algorithm 1).

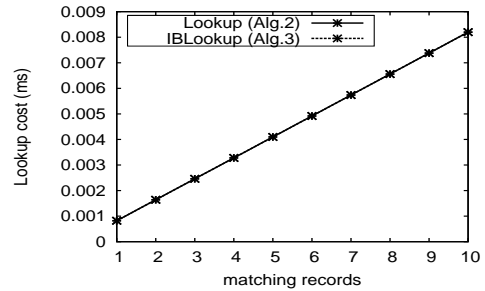


Figure 7: Lookup (Alg. 2) and IBLookup (Alg. 3).

We conclude that, as all operations have linear complexity, our approaches scale efficiently for larger databases and query sets. As a result, one can easily infer results run with even larger parameters, hence, we omit them here.

7.2 First PPSSI Approach vs PIR

We now aim at comparing the efficiency of proposed first PPSSI approach (Section 5) to that of related work – SPIR. Recall that first PPSSI approach provides very similar privacy properties of SPIR. Indeed, both PPSSI and SPIR hide client’s access patterns to the server and also protect privacy of server’s data (with respect to records not matching the queries). However, one possible criticism against our side is that the communication overhead is *linear* in the size of the database size, whereas, SPIR incurs *sub-linear* communication overhead. Remark, however, that: (1) SPIR does not support keyword search, and (2) SPIR introduces a remarkably higher computation overhead, which ends up “overshadowing” the advantage in the communication complexity. To support the latter claim, we compare the overall performance of our first PPSSI approach with that of Gentry and Ramzan’s single-database PIR (GR-PIR) [19], which is, to the best of our knowledge, the most efficient single-database PIR. Specifically, GR-PIR [19], assuming a database with n records, incurs $O(k + d)$ communication complexity (where $k \leq \log n$ and d is the bit-length of each

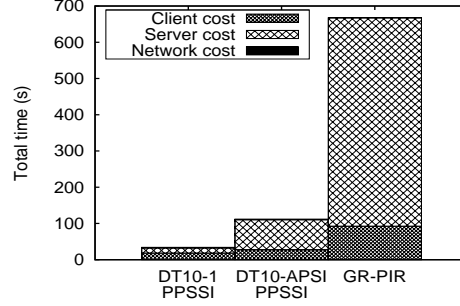


Figure 8: Performance comparison between the first PPSSI approach (Section 5) and GR-PIR [19].

record), and $O(n)$ computation overhead. Also recall that, according to [32], any single-database PIR can be extended to SPIR/OT and we are not aware of any SPIR/OT that is more efficient than GR-PIR.

In our comparison, we use a database with $w = 1024$ records and $m = 5$ attributes. Each record has size $2KB$. We assume the client’s query size is $v = 1024$ and there will be 10 (1%) records matching the query (v_m). On a conservative stance, we choose a relatively slow connection between the client and the server, i.e., a $10Mbps$ link. Remark that we choose 2048-bit modulus and we use RC4 and SHA1.

The result of our comparison is showed in Figure 8 and confirms that our approach is significantly more efficient than GR-PIR. We break down the results into client, server and network transmission cost. Note that, for all schemes, network cost (at the top stack in each bar) is negligible compared to client and server cost. Also observe that GR-PIR imposes a significant overhead on both client and server. We do not show results for larger databases, since: (1) both server and client computational costs will always increase linearly for all schemes, and (2) for very large database, we prefer the approach with the Isolated Box (whose overall performance is evaluated next).

7.3 Second PPSSI approach vs MySQL

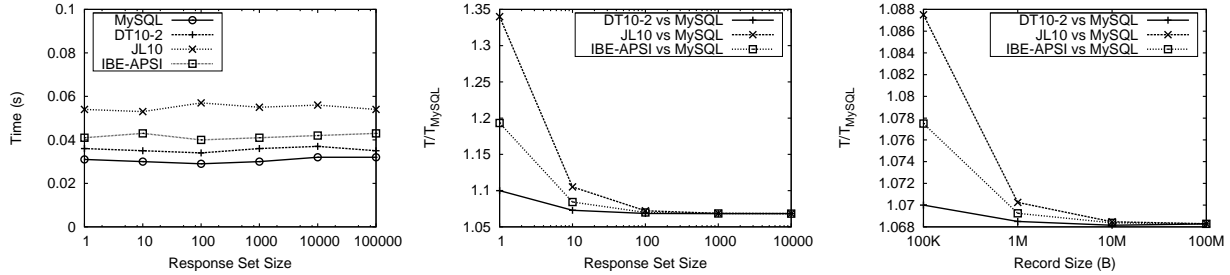
To the best of our knowledge, there is no available approach to PPSSI that combines efficiency with provably secure guarantees and that relies on a non-colluding, untrusted party, such as the Isolated Box. Therefore, we cannot compare our second PPSSI approach for very large databases (Section 6) to any prior work. Nonetheless, we evaluate its performance by measuring it against standard (non privacy-preserving) MySQL.

On a conservative stance, we use MySQL with indexing enabled on each searchable attributes. We run the IB and the server on the same machine. Client is connected to the server and the IB through a $100Mbps$ link. The testing database has 45 searchable attributes and 1 unsearchable attribute (type “LARGEBLOB”) used to pad each record to a uniform size. There are, in total, 100,000 records. All records have the same size, which we vary during experiments.

First, we compare the *index lookup time*, defined as the time between SQL query issuance and the receipt of the first response from the IB. We select a set of SQL queries that return 0, 1, 10, 100, 1000, 10000 ($\pm 10\%$) responses, respectively, and fix each record size at 500KB. Figure 9(a) shows index lookup time for our PPSSI approach (with respect to all underlying (A)PSI-DT instantiations), as well as MySQL, with respect to the response set size. All proposed schemes’ cost are slightly more expensive than MySQL and are independent of the response size.

Next, we test the impact of the response set size on the *total query time*, which we define as the time between SQL query issuance and the arrival of the last response from the IB. Figure 9(b) shows the time for the client to complete a query for a specific response set size divided by the time taken by MySQL (again, with respect to all underlying (A)PSI-DT instantiations). Results gradually converge to 1.1 for increasing response set sizes, i.e., our approach is only 10% slower than standard MySQL. This is because the extra delay incurred by cryptographic operations (in the oblivious evaluation of `TOKEN`) is amortized by subsequent data lookups and decryptions. Note that we can also infer the impact of various client query set size by multiplying the client query set size with each single query delay.

Last, we test the impact of record size on the total query time. We fix response set size at 100 and vary each record size between $100KB$ and $100MB$. Figure 9(c) shows the ratio between our PPSSI approach and MySQL, once more with respect to all underlying (A)PSI-DT instantiations. Again, results gradually converge well below 1.1 with increasing record size. This occurs because, with bigger records, the overhead of record decryption becomes the “bottleneck”.



(a) Index lookup speed comparison. (b) Comparison to MySQL w.r.t. response size. (c) Comparison to MySQL w.r.t. record size.

[DT10-2, JL10, IBE-APSI labels indicate the instantiation used for the Token function in PPSSI]

Figure 9: Performance comparison between the second PPSSI approach (Section 6) and MySQL.

8 Conclusion

In this paper, we proposed secure and efficient techniques for Privacy-Preserving Sharing of Sensitive Information (PPSSI), which enable a client and a server to exchange information without leaking more than the required minimum of information. Privacy guarantees are formally defined and achieved with provable security.

We implemented two variants of PPSSI: one is geared for small/medium-size data sets, while the other minimizes communication overhead, as well as liability issues, for very large databases. The latter introduces a non-colluding, untrusted party – the Isolated Box – which can be implemented as a piece of secure hardware.

Finally, we presented extensive experimental results, which confirmed that our PPSSI approaches are efficient enough to be used in real-world applications. Our future work includes supporting versatile query predicates (e.g., conjunctive queries) as well as fuzzy queries over non-normalized data.

References

- [1] D. Asonov and J.-C. Freytag. Almost optimal private information retrieval. In *Privacy Enhancing Technologies*, 2003.
- [2] M. Belenkiy, J. Camenisch, M. Chase, M. Kohlweiss, A. Lysyanskaya, and H. Shacham. Randomizable proofs and delegatable anonymous credentials. In *Crypto*, 2009.
- [3] D. Boneh, G. Di Crescenzo, R. Ostrovsky, and G. Persiano. Public key encryption with keyword search. In *Eurocrypt*, 2004.
- [4] D. Boneh and M. K. Franklin. Identity-based encryption from the weil pairing. *SIAM Journal of Computing*, 32(3):586–615, 2003.
- [5] Burt Kaliski. TWIRL and RSA Key Size. <http://www.rsa.com/rsalabs/node.asp?id=2004>.
- [6] J. Camenisch and G. M. Zaverucha. Private intersection of certified sets. In *Financial Cryptography'09*, 2009.
- [7] Caslon Analytics. Consumer Data Losses. <http://www.caslon.com.au/datalossnote.htm>.
- [8] B. Chor, N. Gilboa, and M. Naor. Private information retrieval by keywords. *Manuscript*, 1998.
- [9] B. Chor, E. Kushilevitz, O. Goldreich, and M. Sudan. Private information retrieval. *Journal of the ACM*, 45(6):965–981, 1998.
- [10] S. Chow, J. Lee, and L. Subramanian. Two-party computation model for privacy-preserving queries over distributed databases. In *NDSS*, 2009.
- [11] J. Daeman and V. Rijmen. AES proposal: Rijndael. 1999.
- [12] E. De Cristofaro, S. Jarecki, J. Kim, and G. Tsudik. Privacy-preserving policy-based information transfer. In *PETS*, 2009.
- [13] E. De Cristofaro, J. Kim, and G. Tsudik. Linear-Complexity Private Set Intersection Protocols Secure in Malicious Model. In *Asiacrypt*, 2010.
- [14] E. De Cristofaro and G. Tsudik. Practical private set intersection protocols with linear complexity. In *FC*, 2010.
- [15] U. Feige, J. Killian, and M. Naor. A minimal model for secure computation (extended abstract). In *STOC*, 1994.
- [16] M. Fischlin, B. Pinkas, A.-R. Sadeghi, T. Schneider, and I. Visconti. Secure set intersection with untrusted hardware tokens. In *CT-RSA*, 2011.
- [17] M. Freedman, Y. Ishai, B. Pinkas, and O. Reingold. Keyword search and oblivious pseudorandom functions. In *TCC*, 2005.
- [18] M. Freedman, K. Nissim, and B. Pinkas. Efficient private matching and set intersection. In *Eurocrypt*, 2004.
- [19] C. Gentry and Z. Ramzan. Single-database private information retrieval with constant communication rate. In *ICALP*, 2005.

- [20] Y. Gertner, Y. Ishai, E. Kushilevitz, and T. Malkin. Protecting data privacy in private information retrieval schemes. In *STOC*, 1998.
- [21] H. Hacigümüş, B. Iyer, C. Li, and S. Mehrotra. Executing SQL over encrypted data in the database-service-provider model. In *SIGMOD*, 2002.
- [22] C. Hazay and Y. Lindell. Constructions of truly practical secure protocols using standardsmartcards. In *CCS*, 2008.
- [23] C. Hazay and Y. Lindell. Efficient protocols for set intersection and pattern matching with security against malicious and covert adversaries. In *TCC*, 2008.
- [24] B. Hore, S. Mehrotra, and G. Tsudik. A privacy-preserving index for range queries. In *VLDB*, 2004.
- [25] S. Jarecki and X. Liu. Efficient Oblivious Pseudorandom Function with Applications to Adaptive OT and Secure Computation of Set Intersection. In *TCC*, pages 577–594, 2009.
- [26] S. Jarecki and X. Liu. Fast secure computation of set intersection. In *SCN*, 2010.
- [27] M. Kantarcioğlu and C. Clifton. Assuring privacy when big brother is watching. In *DMKD*, 2003.
- [28] M. Kantarcioğlu and J. Vaidya. An architecture for privacy-preserving mining of client information. In *CRPIT*, 2002.
- [29] L. Kissner and D. Song. Privacy-preserving set operations. In *CRYPTO*, 2005.
- [30] M. Murugesan, W. Jiang, C. Clifton, L. Si, and J. Vaidya. Efficient privacy-preserving similar document detection. *VLDB*, 2010.
- [31] S. Nagaraja, P. Mittal, C. Hong, M. Caesar, and N. Borisov. BotGrep: Finding Bots with Structured Graph Analysis. In *Usenix Security*, 2000.
- [32] M. Naor and B. Pinkas. Oblivious Transfer and Polynomial Evaluation. In *STOC*, 1999.
- [33] F. Olumofin and I. Goldberg. Privacy-preserving queries over relational databases. In *PETS*, 2010.
- [34] M. Rabin. How to exchange secrets by oblivious transfer. TR-81, Harvard Aiken Computation Lab, 1981.
- [35] M. Raykova, B. Vo, S. Bellovin, and T. Malkin. Secure anonymous database search. In *CCSW*, 2009.
- [36] R. Rivest. The RC4 Encryption Algorithm. 1992.
- [37] Sherri Davidoff. What Does DHS Know About You?
<http://philosecurity.org/2009/09/07/what-does-dhs-know-about-you>.
- [38] D. Song, D. Wagner, and A. Perrig. Practical techniques for searches on encrypted data. In *S&P*, 2000.

A Details of Token Instantiations based on (A)PSI-DT's

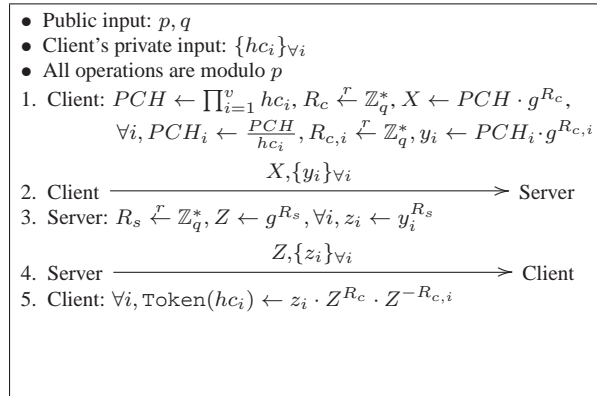


Figure 10: Oblivious computation of $\text{Token}(\cdot)$ using DT10-1.

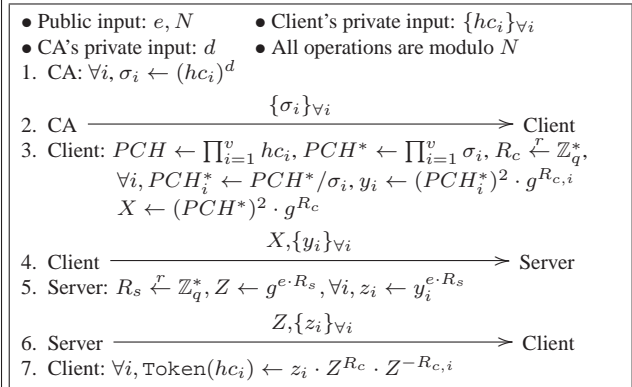


Figure 11: Oblivious computation of $\text{Token}(\cdot)$ using DT10-APSI.

B Security Analysis of First PPSI Approach

We use q_i to denote the i th query of the form $(attr, val)$ issued by the client and use Q_i to denote all records matching query q_i .

<ul style="list-style-type: none"> • Public input: e, N • Client's private input: $\{hc_i\}_{\forall i}$ • Server's private input: $d = e^{-1} \bmod \phi(N)$ where $\phi(N)$ denotes the order of \mathbb{Z}_N^* <ol style="list-style-type: none"> 1. Client: $\forall i, r_i \xleftarrow{r} \mathbb{Z}_N^*, y_i \leftarrow (r_i)^e \cdot hc_i \bmod N$ $\xrightarrow{\{y_i\}_{\forall i}}$ Server 2. Client $\xrightarrow{\{y_i\}_{\forall i}}$ Server 3. Server: $\forall i, z_i \leftarrow (y_i)^d \bmod N$ $\xrightarrow{\{z_i\}_{\forall i}}$ Client 4. Server $\xrightarrow{\{z_i\}_{\forall i}}$ Client 5. Client: $\forall i, \text{Token}(hc_i) \leftarrow z_i \cdot r_i^{-1} \bmod N$ 	<ul style="list-style-type: none"> • Public input: p, q • Client's private input: $\{hc_i\}_{\forall i}$ • Server's private input: $k \in \mathbb{Z}_q^*$ <ol style="list-style-type: none"> 1. Client: $\forall i, \alpha_i \xleftarrow{r} \mathbb{Z}_q^*,$ $y_i \leftarrow ((hc_i)^{(p-1)/q})^{\alpha_i} \bmod p$ $\xrightarrow{\{y_i\}_{\forall i}}$ Server 2. Client $\xrightarrow{\{y_i\}_{\forall i}}$ Server 3. Server: $\forall i, z_i \leftarrow y_i^k \bmod \mathbb{Z}_p^*$ $\pi \leftarrow ZKPK\{k \{z_i = y_i^k\}_{\forall i}\}$ $\xrightarrow{\{z_i\}_{\forall i}, \pi}$ Client 4. Server $\xrightarrow{\{z_i\}_{\forall i}, \pi}$ Client 5. Client: Aborts if π doesn't verify. $\forall i, \text{Token}(hc_i) \leftarrow z_i^{1/\alpha_i} \bmod p$ 	<ul style="list-style-type: none"> • Public input: $P, Q = P^s$ • Client's private input: $\{hc_i\}_{\forall i}$ • CA's private input: s. • Server's private input: z <ol style="list-style-type: none"> 1. CA: $\forall i, \sigma_i \leftarrow (hc_i)^s$ $\xrightarrow{\{\sigma_i\}_{\forall i}}$ Client 2. CA $\xrightarrow{\{\sigma_i\}_{\forall i}}$ Client 3. Server: $R \leftarrow P^z$ (Offline) 4. Server \xrightarrow{R} Client 5. Client: $\forall i, \text{Token}(hc_i) \leftarrow \hat{e}(R, \sigma_i)$
---	--	--

Figure 12: Oblivious computation of $\text{Token}(\cdot)$ using DT10-2.

Figure 13: Oblivious computation of $\text{Token}(\cdot)$ using JL10.

Figure 14: Oblivious computation of $\text{Token}(\cdot)$ using IBE-APSI.

B.1 Security against Honest-but-Curious/Malicious Client

We define against Honest-but-Curious/Malicious client by comparing its view under real model with that under ideal model. In the ideal model, there is a trusted third party (TTP) serving as an honest server who, in response to the query q_i , only replies Q_i .

We first consider Honest-but-Curious adversary and analyze malicious adversary at the end of this section. We define a simulator SIM that attempts to simulate to a real-model client based on output from ideal-model TTP as follows:

Simulator SIM:

SIM is given input $\{q_1, \dots, q_n\}$

1. SIM picks all the secret and public parameters.
2. SIM interacts with \mathcal{A} as a real-model server during oblivious computation of Token (step 1 of Figure 1).
3. SIM sends $\{q_1, \dots, q_n\}$ to the TTP and receives $\{Q_1, \dots, Q_n\}$.
4. SIM runs an arbitrary function on $\{Q_1, \dots, Q_n\}$ and outputs the result to the client.

We then define an experiment for any adversary \mathcal{A} :

The experiment $\text{SPriv}_{\mathcal{C}, \mathcal{A}}$:

1. The adversary \mathcal{A} outputs to the challenger a list of queries $\{q_1, \dots, q_n\}$.
2. The challenger chooses a random bit $b \xleftarrow{r} \{0, 1\}$ and does one of the following:
 - (a) If $b = 0$, then the challenger interacts with \mathcal{A} as a real-model server.
 - (b) If $b = 1$, then the challenger interacts with \mathcal{A} as $\text{SIM}(\{Q_1, \dots, Q_n\})$.
3. The adversary \mathcal{A} outputs a bit b' .
4. The output of the experiment is defined to be 1 if $b' = b$, and 0 otherwise.

Definition 1 *The first PPSSI approach is secure against honest-but-curious client if, for all probabilistic polynomial-time adversaries \mathcal{A} , there exists a probabilistic polynomial-time simulator SIM such that*

$$\Pr[\text{SPriv}_{\mathcal{C}, \mathcal{A}} = 1] \leq \frac{1}{2} + \epsilon$$

This definition ensures that the client in the real model does not get more or different information than the ideal implementation.

Theorem 1 *If the hash function $H(\cdot)$, $H_1(\cdot)$, $H_2(\cdot)$, $H_3(\cdot)$ are collision resistant, Enc is a semantic secure encryption, and $Token$ is unpredictable, then first PPSSI approach is secure against any probabilistic polynomial-time honest-but-curious client.*

Proof: Our goal is to construct a simulator SIM such that \mathcal{A} cannot tell the difference between the view when interacting with SIM and the view when interacting with real-model server. Our SIM is constructed as follows:

1. SIM picks all the secret and public parameters on behalf of a real-model server and publish all public parameters.
2. SIM interacts with \mathcal{A} as a real-model server during oblivious computation of $Token$ (step 1 of Figure 1).
3. SIM queries TTP for $\{q_1, \dots, q_n\}$ and gets back $\{Q_1, \dots, Q_n\}$.
4. Let Q denote $\cup_i Q_i$. SIM generates $w - |Q|$ random records of the same length as any other message in Q . Let DB' denote the concatenation of Q and these random records. Note that $|DB'| = w$.
5. Use Algorithm 1 to encrypt DB' and returns encrypted database EDB' to the client.

We first analyze \mathcal{A} 's view between tags in EDB and tags in EDB' . Note that a tag in $LTable$ is computed as $H_1(\text{Token}(H(attr, val)) || ctr)$. For all $(attr, val)$ pairs not queried by $\{q_1, \dots, q_n\}$, the computed tags should be uniformly random unless (1) there exists j such that $H(q_j) = H(attr, val)$; (2) there exists two pairs $(attr', val')$, $(attr'', val'')$ – such that $H_1(\text{Token}(H(attr', val')) || ctr') = H_1(\text{Token}(H(attr'', val'')) || ctr'')$; (3) \mathcal{A} forges $\text{Token}(H(attr, val))$ for certain $(attr, val)$. All these happen with negligible probability if $H(\cdot)$, $H_1(\cdot)$ are collision resistant and $Token$ is unpredictable.

Next we analyze \mathcal{A} 's view between $(\{ek_{j,l}, eind_{j,l}\}_{1 \leq l \leq m}, er_j)_{1 \leq j \leq w}$ in EDB and those in EDB' . For all $ek, eind, er$ whose corresponding tags do not match $\{q_1, \dots, q_n\}$, they should appear uniformly random to \mathcal{A} unless (1) \mathcal{A} breaks symmetric encryption algorithm; (2) finds collision in $H_2(\cdot)$ or $H_3(\cdot)$; (3) \mathcal{A} can forge $\text{Token}(H(attr, val))$ for certain $(attr, val)$. All these happen with negligible probability if $H_2(\cdot)$, $H_3(\cdot)$ are collision resistant, Enc is semantic secure and $Token$ is unpredictable. \square

In order to consider malicious adversary, we need to change the simulator definition and the experiment. In SIM, there is no input of $\{q_1, \dots, q_n\}$ and, in $SPriv_{C, \mathcal{A}}$, there is no step 1. Note, for the first PPSSI approach, it is secure against malicious adversary only if [13] is used for oblivious computation of $Token$.

Theorem 2 *If oblivious computation of $Token$ protocol is secure against malicious client, the hash function $H(\cdot)$, $H_1(\cdot)$, $H_2(\cdot)$, $H_3(\cdot)$ are collision resistant and Enc is a semantic secure encryption, then first PPSSI approach is secure against any probabilistic polynomial-time malicious client.*

Proof: SIM construction is the same as that in the proof for theorem 1 except that, in step 2, SIM extracts all $\{q_1, \dots, q_n\}$ from the ZKPK sent by \mathcal{A} , which requires rewinding of \mathcal{A} . Then the proof follows that for Theorem 1. \square

B.2 Security against Honest-but-Curious/Malicious Server

Given that the server gets no output from the protocol, the definition of client's privacy requires simply that the server cannot distinguish between cases in which the client has different inputs.

We define an experiment for any adversary \mathcal{A} :

The experiment $SPriv_{S, \mathcal{A}}$:

1. The adversary \mathcal{A} chooses its own database DB and outputs to the challenger two list of queries – (q_1^0, \dots, q_n^0) and (q_1^1, \dots, q_n^1) .
2. The challenger chooses a random bit $b \xleftarrow{r} \{0, 1\}$ and does one of the following:
 - (a) If $b = 0$, then the challenger interacts with \mathcal{A} as a honest client using queries (q_1^0, \dots, q_n^0) .
 - (b) If $b = 1$, then the challenger interacts with \mathcal{A} as a honest client using queries (q_1^1, \dots, q_n^1) .
3. The adversary \mathcal{A} outputs a bit b' .
4. The output of the experiment is defined to be 1 if $b' = b$, and 0 otherwise.

Definition 2 *The first PPSSI approach is secure against honest-but-curious/malicious server if, for all probabilistic polynomial-time adversaries \mathcal{A} ,*

$$\Pr[\text{SPriv}_{S,\mathcal{A}} = 1] \leq \frac{1}{2} + \epsilon$$

Theorem 3 *If oblivious computation of Token function is secure against any probabilistic polynomial-time honest-but-curious/malicious server, the first PPSSI approach is secure against any probabilistic polynomial-time honest-but-curious/malicious server.*

Proof: In the first PPSSI approach, the only messages \mathcal{A} gets from the client is during oblivious Token computation. If oblivious computation of Token function is secure against any probabilistic polynomial-time honest-but-curious/malicious server, the messages \mathcal{A} receives from the client should be perfectly hidden by randomness. Therefore the theorem follows. \square

C Security Analysis of Second PPSSI Approach

Since we do not consider collusion, the security against Honest-but-Curious/Malicious client and server follows exactly from Theorem 1, 2, 3. So we only discuss security against Honest-but-Curious/Malicious Isolated Box.

C.1 Security against Honest-but-Curious/Malicious Isolated Box (IB)

We define security against Honest-but-Curious/Malicious Isolated Box (IB) by comparing its view when interacting with an honest client and an honest server with its view when interacting with a simulator SIM.

Simulator SIM:

SIM is given $|X_U|$ for any $U \subseteq \{0, \dots, n\}$ where $X_U = \cap_{i \in U} Q_i$.

1. SIM outputs an encrypted database \mathbf{EDB}' to \mathcal{A} .
2. SIM interacts with \mathcal{A} as a client, simulating queries $\{q_1, \dots, q_n\}$ (even though SIM does not know $\{q_1, \dots, q_n\}$).

Note, in the above definition, the only information SIM knows is the cardinality of X_U which is defined as the intersection of a subset of query answers.

We then define an experiment for any adversary \mathcal{A} :

The experiment $\text{SPriv}_{\text{IB},\mathcal{A}}$:

1. The adversary \mathcal{A} outputs to the challenger a database DB and a list of queries $\{q_1, \dots, q_n\}$.
2. The challenger chooses a random bit $b \xleftarrow{r} \{0, 1\}$ and does one of the following:
 - (a) If $b = 0$, then the challenger interacts with \mathcal{A} as an honest client and an honest server.
 - (b) If $b = 1$, then the challenger computes $\{Q_1, \dots, Q_n\}$ based on DB , derives all intersections X_U for all $U \subseteq \{1, \dots, n\}$ and interacts with \mathcal{A} as $\text{SIM}(\{|X_U|\}_{\forall U \subseteq \{1, \dots, n\}})$.
3. The adversary \mathcal{A} outputs a bit b' .
4. The output of the experiment is defined to be 1 if $b' = b$, and 0 otherwise.

Definition 3 *The second PPSSI approach is secure against honest-but-curious/malicious IB if, for all probabilistic polynomial-time adversaries \mathcal{A} , there exists a probabilistic polynomial-time simulator SIM such that*

$$\Pr[\text{SPriv}_{\text{IB},\mathcal{A}} = 1] \leq \frac{1}{2} + \epsilon$$

Theorem 4 *If the hash function $H(\cdot)$, $H_1(\cdot)$, $H_2(\cdot)$, $H_3(\cdot)$ are collision resistant, Enc is a semantic secure encryption, and $Token$ is unpredictable, then the second PPSSI approach is secure against any probabilistic polynomial-time honest-but-curious/malicious IB.*

Proof: Our goal is to construct a simulator SIM such that \mathcal{A} cannot tell the difference between the view when interacting with SIM and the view when interacting with an honest client and an honest server. Our SIM is constructed as follows:

1. SIM creates **EDB'**:

- Pick w random messages of same length as encrypted messages.
- Then create **LTable'** = $\{(tag'_{j,l}, ek'_{j,l}, eind'_{j,l})\}_{1 \leq j \leq w, 1 \leq l \leq m}$ where $tag'_{j,l} \in_R \{0, 1\}^{l_h}$, $ek'_{j,l} \in_R \{0, 1\}^{l_e}$, $eind'_{j,l} \in_R \{0, 1\}^{l_d}$, l_h is the output length of hash function, l_e is the output length of encryption function.

2. For each query q_i , SIM prepares the matching tag set $T_i = \{tag_1^i \dots tag_{|Q_i|}^i\}$ such that $|\cap_{i \in U} T_i| = |X_U|$ for any $U \subseteq \{0, \dots, n\}$ as follows:

- For all U , compute $|\hat{X}_U|$ where $\hat{X}_U = X_U \setminus \cup_{|U'| > |U|} X_{U'}$. Given $|X_U|$, $|\hat{X}_U|$ can be computed as

$$\begin{aligned} |\hat{X}_U| &= |X_U| - |X_U \cap (\cup_{|U'| > |U|} X_{U'})| \\ &= |X_U| - \left(\sum_{|U'| > |U|} |X_U \cap X_{U'}| - \sum_{|U'_1| > |U|, |U'_2| > |U|, U'_1 \neq U'_2} |(X_U \cap X_{U'_1}) \cap (X_U \cap X_{U'_2})| \right. \\ &\quad \left. + \dots + (-1)^{\binom{n}{n} + \dots + \binom{n}{|U|+1}} \cdot |\cap_{|U'| > |U|} (X_U \cap X_{U'})| \right) \end{aligned}$$

where $|X_{U_1} \cap \dots \cap X_{U_n}| = |X_{U_1 \cup \dots \cup U_n}|$. It is easy to observe that $\sum_{U \subseteq \{1, \dots, n\}} |\hat{X}_U| = |\cup_{j=1}^n Q_j|$

- Randomly pick $\sum_{\forall U} |\hat{X}_U|$ different tags from **LTable'** and store them in Y . For each U , initialize \hat{Q}_U as follows:
 - (a) Pick $|\hat{X}_U|$ distinct tags from Y and add them to \hat{Q}_U .
 - (b) Update $Y \leftarrow Y \setminus \hat{Q}_U$.
- For $\lambda = 1, \dots, n$, set $T_\lambda = \cup_{\lambda \in U} \hat{Q}_U$ and append a random tag (used to terminate a query) which is different from all tags in **LTable'** to T_λ . Note $|T_\lambda| = |Q_\lambda| + 1$.

3. SIM plays the role of a client as follows: for the λ th query, make $|T_\lambda|$ probes where θ th probe is the θ th element in T_λ .

We first analyze the view of \mathcal{A} between tags in **EDB** and those in **EDB'**. The distribution of tags in **EDB** and those in **EDB'** is the same unless one of the following happens: (1) there exists $(attr_i, val_i) \neq (attr_j, val_j)$ but $H(attr_i, val_i) = H(attr_j, val_j)$; (2) $H(attr_i, val_i) \neq H(attr_j, val_j)$ but $H_1(\text{Token}(H(attr', val')) || ctr') = H_1(\text{Token}(H(attr'', val'')) || ctr'')$; (3) \mathcal{A} forges $\text{Token}(H(attr, val))$ for certain $(attr, val)$. All these happen with negligible probability if $H(\cdot)$, $H_1(\cdot)$ are collision resistant and Token is unpredictable.

Next we analyze \mathcal{A} 's view between $(\{ek_{j,l}, eind_{j,l}\}_{1 \leq l \leq m}, er_j)_{1 \leq j \leq w}$ in **EDB** and those in **EDB'**. They should appear uniformly random to \mathcal{A} unless (1) \mathcal{A} breaks symmetric encryption algorithm; (2) \mathcal{A} finds collision in $H_2(\cdot)$ or $H_3(\cdot)$ (which breaks one-time-pad encryption); (3) \mathcal{A} can forge $\text{Token}(H(attr, val))$ for certain $(attr, val)$. All these happen with negligible probability if $H_2(\cdot)$, $H_3(\cdot)$ are collision resistant, Enc is semantic secure and Token is unpredictable.

Last we show that \mathcal{A} cannot distinguish the way that an honest client's queries are answered using **EDB** and the way that SIM's queries are answered using **EDB'**. For an honest client's query q_i , there are $|Q_i|$ matches in **EDB**. For the SIM's i th query, it makes $|T_i|$ probes and there will be $|T_i| - 1$ matches. Since $|T_i| - 1 = |Q_i|$ and \mathcal{A} cannot distinguish er_j from er'_j , the view in the real protocol and that in the interaction with SIM are identical. \square

D Comparison of State-of-the-art PSI-DT

In the following, we review and compare state-of-the-art PSI protocols and focus on PSI-DT variants. We assume client and server set sizes are v and w , respectively.

D.1 PSI-DT without Pre-distribution

FNP04. Freedman, Nissim, and Pinkas [18] use *oblivious polynomial evaluation* to implement PSI. Their approach can be slightly modified to support PSI-DT. The modified protocol – denoted as FNP04 – works as follows: the client first setups an additively homomorphic encryption scheme, such as Paillier, with key pair (pk_c, sk_c) . Client defines a polynomial $f(y) = \prod_{i=1}^v (y - c_i) = \sum_{i=0}^v a_i y^i$ whose roots are its inputs. It encrypts each coefficient a_i under its public key pk_c and sends encrypted coefficients $\{Enc_{pk_c}(a_i)\}_{i=0}^k$ to the server. Since the encryption is homomorphic, the server can evaluate $Enc(f(s_j))$ for each $s_j \in \mathcal{S}$ independently from the client. Then, the server returns $\{(Enc(r_j \cdot f(s_j) + s_j), Enc(r'_j \cdot f(s_j) + data_j))\}_{j=0}^n$ to the client where r_j and r'_j are fresh random numbers for each input in \mathcal{S} . Client, for each returned pair (e_l, e_r) , decrypts e_l by computing $c' = Dec_{sk_c}(e_l)$. Then if $c' \in \mathcal{C}$, the client continues to decrypt e_r and gets the associated data. Otherwise, the client only gets some random value and moves onto the next returned pair. In order to speed up the performance, FNP04 can use modified ElGamal encryption instead of Paillier. Specifically, the client uses g^{a_i} instead of a_i as the input to the ElGamal encryption where g is a generator with order q modulo p . And when it decrypts e_l , it recovers $g^{c'}$. Client can still decide whether $c' \in \mathcal{C}$ by comparing $g^{c'}$ to $g^{c_i}, \forall c_i \in \mathcal{C}$. In terms of data, the server can choose a random key g^{k_j} and uses it to symmetrically encrypt $data_j$. Then the server sends $\{(Enc(r_j \cdot f(s_j) + s_j), Enc(r'_j \cdot f(s_j) + k_j), Enc_{g^{k_j}}(data_j))\}_{j=0}^w$ to the client. If the client can recover g^{k_j} , it can also decrypt $data_j$. Using balanced bucket allocation to speed up operations, client overhead is dominated by $O(v + w) |q|$ -bit mod p exponentiations (in ElGamal). Whereas, server overhead is dominated by $O(w \log \log v) |q|$ -bit mod p exponentiations.

KS05. Kissner and Song [29] also use oblivious polynomial evaluation to construct a variety of set operations. However, their solution is designed for mutual intersection over *multi-set* that may contain duplicate elements, and it is unclear how to adapt it to transfer associated data. Also, their technique incurs quadratic ($O(vw)$) computation (but linear communication) overhead. As we use a different method to handle multi-sets (see Section 5) and we only consider one-way PSI, we do not consider KS05 any further.

DT10-1. De Cristofaro and Tsudik present an unlinkable PSI-DT protocol (Fig. 3 in [14]) with linear computation and communication complexities. This protocol, denoted as DT10-1, operates as follows: The setup phase yields primes p (e.g. 1024 bits) and q (e.g. 160 bits), s.t. $q|p-1$, and a generator g with order q modulo p . In the following, we assume computation is done mod p . First, the client sends to the server $X = [(\prod_{i=1}^v H(c_i)) \cdot g^{R_c}]$ where R_c is randomly selected from \mathbb{Z}_q . Also, for each $1 \leq i \leq v$, the client sends $y_i = [(\prod_{l \neq i} H(c_l)) \cdot g^{R_{c:i}}]$, where the $R_{c:i}$'s are random in \mathbb{Z}_q . The server picks a random R_s in \mathbb{Z}_q and replies with $Z = g^{R_s}$ and $y'_i = y_i^{R_s}$ (for every y_i it received). Also, for each item s_j ($1 \leq j \leq w$), it computes $K_{s:j} = (X/H(s_j))^{R_s}$, and sends the tag $t_j = H_1(K_{s:j})$ with the associated data record encrypted under $k_j = H_2(K_{s:j})$. The client, for each of its elements, computes $K_{c:i} = y'_i \cdot Z^{R_c} \cdot Z^{-R_{c:i}}$ and the tag $t'_i = H_1(K_{c:i})$. Only if c_i is in the intersection (i.e., there exists an element $s_j = c_i$), the client finds a pair of matching tags (t'_i, t_j) . Besides learning the elements intersection, the client can decrypt associated data records by key $H_2(K_{c:i})$. Client overhead amounts to $O(v) |q|$ -bit modulo p exponentiations and multiplications and server overhead is $O(v + w) |q|$ -bit modulo p exponentiations.

D.2 PSI-DT with Pre-distribution

JL09. Jarecki and Liu [25] (following the idea in [23]) give a PSI-DT based on Oblivious PRF (OPRF) [17]. We denote this protocol as JL09 (and present the improved OPRF construction discussed in [2]). Recall that an OPRF is a two-party protocol that securely computes a pseudorandom function $f_k(\cdot)$, on key k contributed by a server and input x contributed by a client, such that the server learns nothing about x , while the client learns $f_k(x)$. The main idea is the following: For every item $s_j \in \mathcal{S}$, the server publishes a set of pair $\{H_1(f_k(s_j)), Enc_{H_2(f_k(s_j))}(data_j)\}$. Then, the client, for every item $c_i \in \mathcal{C}$, obtains $f_k(c_i)$ by OPRF with the server. As a result, the client can use $H_1(f_k(c_i))$ to check if $c_i \in \mathcal{C} \cap \mathcal{S}$ and if so then it uses $H_2(f_k(c_i))$ to recover $data_j$. JL09 incurs $O(w + v)$ server exponentiations, and $O(v)$ client exponentiations. Exponentiations are $|N|$ -bit modulo N^2 , where N is the RSA modulus.

JL10. Another recent work by Jarecki and Liu [26] (denoted as JL10) leverages an idea similar to JL09 [25] to achieve PSI-DT. Instead of using OPRF, JL10 uses the newly-introduced *Parallel Oblivious Unpredictable Function* (POUF), $f_k(x) = (H(x)^k \bmod p)$, in the Random Oracle Model. In order to obliviously compute $f_k(x)$, the client first picks a random exponent α and sends $y_j = H(c_j)^\alpha$ to the server. The server replies to the client with $z_j = (y_j)^k$. Then the client recovers $f_k(x) = z^{1/\alpha}$. The computational complexity of this protocol amounts to $O(v)$ online exponentiations for both server and client, as the server can pre-process (offline) its $O(w)$ exponentiations. Exponentiations are q -bit

modulo p , similar to DT10-1.

DT10-2. In Fig. 4 of [14], De Cristofaro and Tsudik present a PSI-DT based on blind-RSA signatures in the Random Oracle Model (ROM). We denote this protocol as DT10-2. The protocol uses the hash of RSA signatures as a PRF in ROM and achieves the same asymptotic complexities as DT10-2 and JL10, but (1) the server now computes RSA signatures (e.g., 1024-bit exponentiations), and (2) client workload is reduced to only multiplications if the RSA public key, e , is chosen short enough (e.g., $e = 3$).

In summary, we consider JL09, JL10 and DT10-2 in the context of PSI-DT *with* pre-distribution. Note that, although faster than protocols without pre-distribution, these protocols do not achieve Server Unlinkability.

D.3 APSI-DT without Pre-distribution

DT10-APSI. In Fig.2 of [14], De Cristofaro and Tsudik also present an APSI-DT technique mirroring its PSI-DT counterpart, DT10-1. We denote this protocol as DT10-APSI. It operates as follows: the client first obtains authorization from the court for its element c_i , where an authorization corresponds to an RSA-signature: $\sigma_i = H(c_i)^d$. Then, the client sends the server $X = [(\prod_{i=1}^v \sigma_i) \cdot g^{R_c}]$ for a random R_c . Then, for each element c_i , it sends $y_i = [(\prod_{l \neq i} \sigma_l) \cdot g^{R_{c:i}}]$, where the $R_{c:i}$'s are additional random values. The server picks a random value, R_s , and replies with $Z = g^{eR_s}$, $y'_i = y_i^{eR_s}$ (for each received y_i). Also, for each element s_j , she computes $K_{s:j} = (X^e / H(s_j))^{R_s}$, and sends the *tag* $t_j = H_1(K_{s:j})$ and the associated data record encrypted under the key $k_j = H_2(K_{s:j})$. Client, for each of its elements, computes $K_{c:i} = y'_i \cdot Z^{R_c} \cdot Z^{-R_{c:i}}$ and the tag $t'_i = H_1(K_{c:i})$. Client can find a pair of matching tag (t'_i, t_j) only if c_i is in the intersection and σ_i is a valid signature on c_i . Besides learning the elements in the intersection, the client can decrypt associated data records. The computation overhead is $O(v)$ exponentiations for the client, and $O(v + w)$ – for the server. Exponentiations are $|N|$ -bit modulo N , where N is the RSA modulus.

CZ09. Camenisch and Zaverucha [6] provide mutual set intersection with authorization on both parties' input. The proposed protocol builds upon oblivious polynomial evaluation and has quadratic computation and communication overhead. Also, it does not provide data transfer.

As a result, we only consider the DT10-APSI protocol in the context of APSI-DT *without* pre-distribution. Note that DT10-APSI provides both Server and Client Unlinkability, as well as Forward Security.

D.4 APSI-DT with Pre-distribution

IBE-APSI. The protocol in Fig. 5 of [12] presents a protocol based on Boneh-Franklin Identity-based Encryption [4], which can be adapted to APSI-DT with pre-distribution. We denote this protocol as IBE-APSI. Note that such a construct is described in the context of a different primitive – Privacy-Preserving Information Transfer (PPIT). However, it can be converted to APSI-DT.

First, the authorization authority (acting as the IBE PKG) generates a prime q , two groups $\mathbb{G}_1, \mathbb{G}_2$ of order q , a bilinear map $e : \mathbb{G}_1 \times \mathbb{G}_1 \rightarrow \mathbb{G}_2$. A random $s \in \mathbb{Z}_q$ is selected as a secret master key. Then, a random generator $P \in \mathbb{G}_1$ is chosen, and Q is set such that $Q = s \cdot P$. (P, Q) are public parameters. Client obtains authorization for an element c_i as an IBE secret key, $\sigma_i = s \cdot H(c_i)$. In the pre-distribution phase, the server first selects a random $z \in \mathbb{G}_1$ and then, for each $(s_j, data_j)$, publishes (t_j, e_j) where $t_j = H_1(e(Q, H(s_j))^z)$ and e_j is the IBE encryption of $data_j$ under identifier s_j . Then, the server gives the client $R = zP$ and the client computes $t'_i = H_1(e(R, \sigma_i))$. For any t'_i , s.t. $t'_i = t_j$, the client can decrypt e_j . The protocol can be speeded up by encrypting e_j under symmetric key $H_2(e(Q, H(s_j))^z)$. The computation overhead for the client amounts to $O(v)$ pairing operations, while there is no online overhead for the server.

Remark that IBE-APSI has two drawbacks compared to APSI-DT: it provides neither Server Unlinkability nor Forward Security.

D.5 Benchmark of (A)PSI-DTs

In this section, we benchmark several (A)PSI-DT protocols and compare their performance through experimental results. During the process, we try to identify the most efficient (A)PSI-DT protocols (with or without pre-distribution), and select the building blocks of our PSSI solutions.

Candidate Protocols. We discuss efficient implementation of the following (A)PSI-DT protocols:

	w/o Pre-Distribution	w/ Pre-Distribution
PSI-DT	FNP04 ([18]), DT10-1 (Fig.3 in [14])	JL09 ([25]), JL10 ([26]), DT10-2 (Fig.4 in [14])
APSI-DT	DT10-APSI (Fig.2 in [14])	IBE-APSI (Fig.5 in [12])

Table 7: Candidate PSI-DT and APSI-DT protocols.

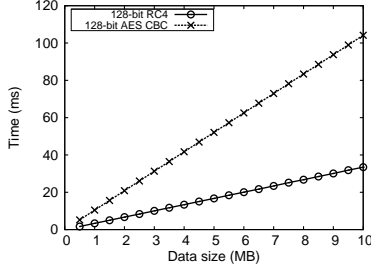


Figure 15: Symmetric key en-/de-cryption performance.

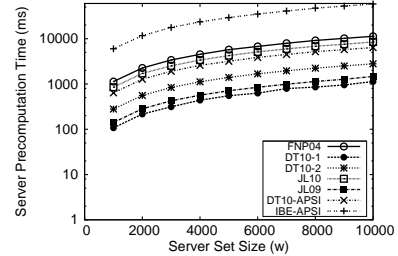


Figure 16: Server pre-computation overhead.

Each protocol was implemented in C++ using GMP (ver. 5.01) and PBC (ver. 0.57) libraries. All benchmarks were collected on a Ubuntu 9.10 desktop platform with Intel Xeon E5420 CPU (2.5GHz and 6MB cache) and 8GB RAM.

For protocols supporting data transfer, data associated with each server element can be arbitrarily long. Also, performance of some protocols is dominated by each element’s data size, rather than set size (e.g., in FNP04). In order to obtain a fair comparison, we need to capture the “intrinsic” cost of each protocol. To this end, we employ the following strategy to eliminate data size effects: First, in all protocols, we encrypt each element’s data with a distinct random symmetric key and consider these keys as the new associated data. Assuming that a different key is selected at each interaction, this technique does not violate Server Unlinkability. This way, the computation cost of each protocol is measured based on the same fixed-length key, regardless of data size. In our experiments, we set symmetric key size to 128 bits.

As a result, each protocol execution involves additional overhead of symmetric en-/de-cryption of records. Figure 15 compares the resulting overhead (for variable data sizes), using either RC4 [36] or AES-CBC [11] (with 128-bit keys). Therefore, to estimate the total cost of a protocol, one needs to combine: (1) symmetric encryption overhead, (2) computation cost of each protocol, and (3) data transfer delay for transmitting the encrypted data and PSI values.

We further assume that the client does not perform any pre-computation, while the server performs as much pre-computation on its input as possible. This reflects the reality where client input is (usually) determined in real time, while server input is pre-determined. Figure 16 shows the pre-computation overhead for each protocol.

Next, we evaluate online computation overhead. Figures 17 and 18 present client online computation overhead with respect to client and server input sizes, respectively. Figures 19 and 20 show server online computation overhead with respect to client and server input size, respectively.

Furthermore, Figures 21 and 22 evaluate protocol bandwidth complexity with respect to client and server input sizes. For protocols with pre-distribution, bandwidth consumption (since the transfer of database encryption is performed offline) does not include pre-distribution overhead. Note that, in these figures, we sometimes use the same marker for different protocols to indicate that these protocols share the same value. Client input size v (resp., server input size w) is fixed at 5,000 in figures where x-axis refers to the server (resp., the client) input size.

Finally, note that, in all experiments, we use a 1024-bit RSA modulus and a 1024-bit cyclic-group modulus with a 160-bit subgroup order. All test results are averaged over 10 independent runs. All protocols are instantiated under the assumption of *Honest-but-Curious* (HbC) adversaries and in the *Random Oracle Model* (ROM).

PSI-DT without pre-distribution. We now focus on the comparison between FNP04 and DT10-1. Figures 17-22 show that that FNP04 is much costlier than DT10-1 in terms of client and server online computation as well as bandwidth consumption. For each client set size, DT10-1 client overhead ranges from $460ms$ to $4,400ms$, while FNP04 server overhead – between $1,300ms$ and $15,000ms$. For each chosen server set size, server overhead in DT10-1 is under $1,300ms$, while in FNP04 it exceeds $15,000ms$.

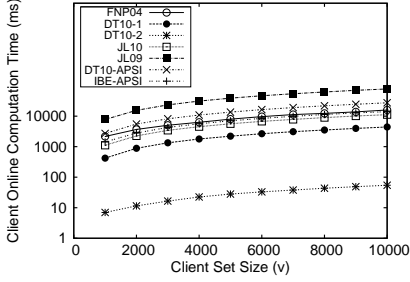


Figure 17: Client online computation w.r.t. client set size.

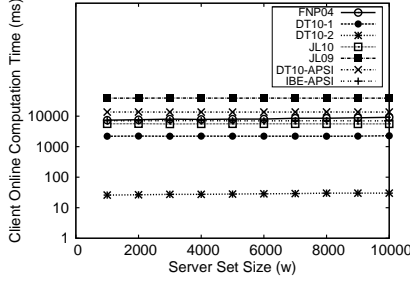


Figure 18: Client online computation w.r.t. server set size.

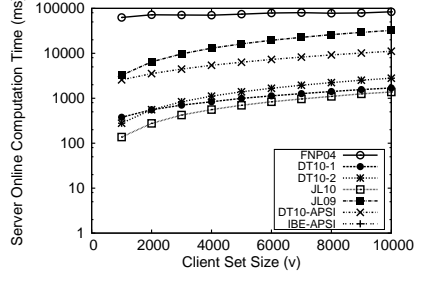


Figure 19: Server online computation w.r.t. client set size.

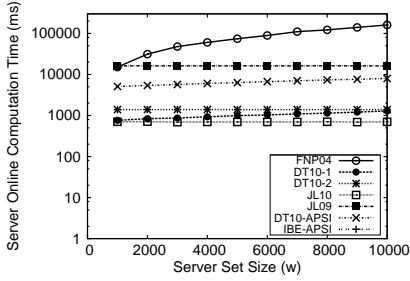


Figure 20: Server online computation w.r.t. server set size.

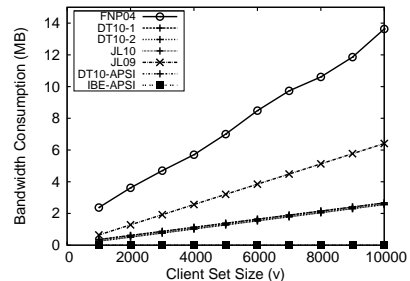


Figure 21: Bandwidth consumption w.r.t. client set size.

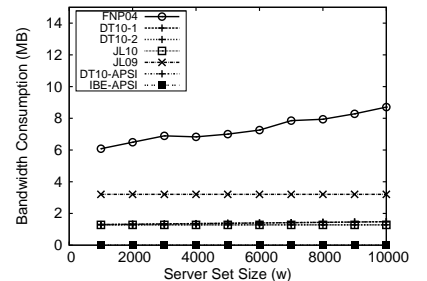


Figure 22: Bandwidth consumption w.r.t. server set size.

PSI-DT with pre-distribution. Next, we compare JL09, JL10 and DT10-2, i.e., PSI-DTs with pre-distribution. Recall that all protocols are instantiated in the HbC model, thus ZKPK's are not included for JL09 and JL10. Figures 17-22 show that DT10-2 incurs client overhead almost two orders of magnitude lower than JL09 and JL10. Indeed, DT10-2 involves two client multiplications for each item, while JL09 performs two heavy homomorphic operations and JL10 – two exponentiations. In JL10, the server online computation overhead results from v 160-bit exponentiations, whereas, in DT10-2, it results from v RSA exponentiations. Since these exponentiations can be speeded up using the Chinese Remainder Theorem, the gap (for server computation overhead) between JL10 and DT10-2 is only double. Summing up server and client computation overhead, DT10-2 results to be the most efficient. In terms of bandwidth consumption, DT10-2 and JL10 are almost the same, while JL09 is slightly more expensive.

APSI-DT without pre-distribution. The only protocol available in this context is DT10-APSI (as discussed in Appendix D.3). Figure 17-20 illustrates that client overhead is determined only by client set size, whereas, server overhead is determined by both client and server set sizes. Note that measurements obtained for APSI-DT naturally mirror those of DT10-1, as the former simply adds authorization of client inputs (by merging signatures into the protocol).

APSI-DT with pre-distribution. The only protocol we evaluate for APSI-DT with data pre-distribution is IBE-APSI (as discussed in Appendix D.4). Figure 17-18 shows that client overhead increases linearly with client set size and does not depend on server set size. Recall that, in IBE-APSI, the server needs to compute pairing operations for each item, independent of client input. Moreover, since these operations can be pre-computed, server-side overhead and bandwidth consumption are negligible, as shown in Figures 19-22.⁴

During the pre-computation phase, the server needs to compute w pairing and exponentiations, which makes pre-computation relatively expensive. Thus, note that, If Server Unlinkability is desired, server would need to repeat, for every interaction, the operations otherwise performed only during pre-computation.

⁴In these figures, y-values for IBE-APSI are all 0 which is out of the scope of the y-axis.