

# Verifiable Delegated Set Intersection Operations on Outsourced Encrypted Data

Qingji Zheng    Shouhuai Xu

Department of Computer Science, University of Texas at San Antonio, TX, USA

**Abstract**—We initiate the study of the following problem: Suppose Alice and Bob would like to outsource their encrypted private data sets to the cloud, and they also want to conduct the set intersection operation on their plaintext data sets. The straightforward solution for them is to download their outsourced ciphertexts, decrypt the ciphertexts locally, and then execute a commodity two-party set intersection protocol. Unfortunately, this solution is not practical.

We therefore motivate and introduce the novel notion of *Verifiable Delegated Set Intersection on outsourced encrypted data* (VDSI). The basic idea is to delegate the set intersection operation to the cloud, while (i) not giving the decryption capability to the cloud, and (ii) being able to hold the misbehaving cloud accountable. We formalize security properties of VDSI and present a construction. In our solution, the computational and communication costs on the users are linear to the size of the intersection set, meaning that the efficiency is optimal up to a constant factor.

**Index Terms**—verifiable set intersection; outsourced encrypted data; verifiable outsourced computing

## I. INTRODUCTION

Cloud computing allows users to outsource their data to the cloud, but the data privacy issue often makes them reluctant to do so. It is therefore natural to encrypt the outsourced data and delegate the heavy-duty computational tasks on the outsourced encrypted data to the cloud. This leads to a general question: How can the cloud execute the delegated functions on outsourced encrypted data, without being given the decryption capability? Although Fully Homomorphic Encryption (FHE) [1]–[3] is promising to tackle this problem, it is not practical enough for applications that involve a large volume of data [4]. Moreover, FHE in general does not solve another important problem: How can we force the cloud to execute the delegated computational functions honestly? This calls for solutions that can hold the misbehaving cloud accountable.

In this paper, we consider the problem of Verifiable Delegated Set Intersection on outsourced encrypted data (VDSI), which can be seen as the cloud version of the well investigated problem of Private Set Intersection (PSI) [5]–[7]. In the setting of PSI, two parties jointly compute the intersection of their private data sets such that they learn the intersection set but nothing else (the sizes of their private data sets may or may not be deemed as confidential [8]).

In the setting of VDSI, two cloud users, Alice and Bob, outsource their encrypted private data sets to the cloud. They would like to conduct the set intersection operation on their plaintext data sets. The straightforward solution would be for them to download their outsourced ciphertexts, decrypt the ciphertexts locally, and then execute a commodity two-party

set intersection protocol. The straightforward solution is not practical, especially when the outsourced data sets are large and when they use wireless systems such as smartphones. Another drawback of this solution is that both Alice and Bob must participate simultaneously. For these reasons, Alice and Bob would prefer delegating the set intersection operation to the cloud, while being able to hold the misbehaving cloud accountable. Note that it is realistic to assume that the cloud is untrusted because it has the incentive not to honestly execute the protocols (e.g., for saving resources or shortening service response time). Moreover, the cloud may have been compromised and the attacker may return Alice and Bob with misleading results.

## A. Our Contribution

We initiate the investigation of a novel notion called VDSI, a useful primitive for delegating the set intersection operation on outsourced ciphertexts to the untrusted cloud. In contrast to the straightforward solution mentioned above, VDSI solves the problem by enabling the cloud to compute the set intersection, but *without* giving the decryption capability to the cloud. As such, VDSI can be seen as a special-purpose homomorphic cryptographic system for use in cloud computing. Since the cloud is untrusted and possibly malicious, VDSI allows Alice and Bob to verify whether the cloud has faithfully computed the delegated set intersection protocol or not.

Specifically, we formally define security properties of VDSI, and present a concrete VDSI scheme. The scheme is based on two ideas: (i) using proxy re-encryption to enable the cloud to compare equality of plaintexts corresponding to two ciphertexts that are encrypted using different public keys; (ii) using a novel variant of cryptographic accumulator, which can be used to verify the membership of *multiple* elements through a single examination and may be of independent value, to allow the cloud to show the correctness of the resulting intersection set.

Our VDSI scheme has two appealing features. First, it does not require the participation of Alice and Bob, because the cloud conducts the delegated computing. Second, it is much more efficient than the straightforward solution mentioned above. Suppose Alice's (Bob's) data set has  $n$  ( $m$ ) elements, and the intersection set has  $k$  elements. Our solution only incurs  $O(k)$  computational and communication costs on Alice and Bob, for decrypting and verifying the results received from the cloud. This means that our solution is optimal (up to a constant factor). In contrast, the straightforward solution incurs  $O(m + n)$  computational and communication costs on

Alice and Bob. Note that it is possible that  $m + n \gg k$ . Experimental evaluation confirms that the VDSI scheme is practical.

We believe that the novel concept of VDSI will inspire many fruitful studies. For example, our solution only achieves a “weak” version of the *function output secrecy* property, which allows the untrusted cloud to launch a *plaintext guessing* attack against Alice’s and Bob’s private data (i.e., the success probability depends on the size of the plaintext space). It is an outstanding open problem to settle down whether or not this weak guarantee is inherent to the problem that VDSI aims to solve; if not, we need to design a better solution that is immune to this attack. Another outstanding open problem is to enforce fine-grained access control over the delegated set intersection operation, which may or may not need to be traded from the verifiability.

### B. Related Work

To the best of our knowledge, this is the first work that considers the PSI problem in the cloud computing setting, where cloud users not only outsource their private data but also outsource their set intersection operations, while being able to hold the dishonest cloud vendors accountable for not faithfully executing the delegated operations. Nevertheless, there are prior studies on related problems.

**Private Set Intersection.** The PSI (private set intersection) problem was initiated in [9] and has become an essential building-block for many applications. Many variants of PSI [6], [8], [10]–[20] have been proposed, with various features (e.g., preventing a malicious party from choosing arbitrary inputs [6], [11], hiding the sizes of the inputs [8], verifiable set operation without preserving data privacy [13], [16]). There have been schemes that aim to reduce the computational and communication complexities (e.g., the RSA-OPRF-based protocol [21], the garbled circuit protocol [20], and the garbled bloom filter protocol [19]). Among the state-of-the-art PSI solutions, the most efficient PSI protocol incurs  $O(m + n)$  computational and communication complexities, where  $m$  and  $n$  are sizes of the respective data sets [19]. We note that [22]–[24] considered the problem of server-aided private set intersection, where cloud users share some secrets with each other to preprocess data sets at the time of outsourcing their data. Such collaborative preprocessing is not needed in the setting of VDSI. Finally, a recent work [25] studies verifiable complex set operations over outsourced plaintext data sets (i.e., the outsourced data is not encrypted). In contrast, we consider outsourced computing on outsourced encrypted data.

**Public Key Encryption with Equality Test.** The problem of public key encryption with equality test is to decide whether two ciphertexts that are encrypted using two different public keys correspond to the same plaintext or not [26]–[29]. In order to enforce access control over the equality test operation, a variant of the problem is to allow the data owners to authorize who can perform the equality test on the outsourced encrypted data [30]. These protocols do not consider the requirement of verifiability on the equality test results, which is crucial to VDSI in the present paper.

**Verifiable Computation.** How to securely and efficiently delegate the computation of a function to a remote server has been under active research [31]–[40]. In these solutions, the data owner pre-processes the inputs to the delegated function in question before outsourcing the data to the cloud, and the cloud needs to prove the correctness of the outcome of a function execution. However, these solutions do not solve the problem studied in this paper because (i) the input to their functions is from a single source and known to the delegator in advance, and (ii) some solutions do not consider privacy of the input. In contrast, our model has the following characteristics: the inputs to the delegated functions include other data owners’ private data sets, which are not known to the delegators in advance. Finally, [41] considered the notion of verifiable private multi-party computation, but not in the setting of outsourcing data and functions to the cloud.

**Paper Organization.** Section II reviews some cryptographic preliminaries. Section III formulates the problem of VDSI and its security properties. Section IV introduces the extended accumulator scheme, which is used as a building-block and may be of independent value. Section V presents a VDSI scheme and analyzes its security, while Section VI evaluates its performance. Section VII concludes the paper.

## II. CRYPTOGRAPHIC PRELIMINARIES

Let  $(e, g, G, G_T, p) \leftarrow \text{MapGen}(1^\ell)$  denote that the bootstrapping algorithm  $\text{MapGen}$  generates a bilinear map  $e : G \times G \rightarrow G_T$ , where  $G$  and  $G_T$  are cyclic groups of order  $p$  which is an  $\ell$ -bit prime,  $g$  is a generator of  $G$ , and the bilinear map  $e$  satisfies (i) for  $a, b \in \mathbb{Z}_p$ ,  $e(g^a, g^b) = e(g, g)^{ab}$ , (ii)  $e(g, g)$  is non-degenerate, and (iii)  $e$  can be efficiently computed. The bilinear map  $e$  is one-way, i.e. the probability of a probabilistic polynomial algorithm inverting  $e$  is negligible, which holds when  $G$  and  $G_T$  are instantiated with Weil or Tate pairing over MNT curves [42]. Table I summarizes the notions for the algorithms and parameters in the VDSI scheme, multi-accumulator scheme and the signature scheme  $\text{Sig}$ .

**Bilinear  $q$ -strong Diffie-Hellman assumption ( $q$ -SDH) [42].** For given  $(e, g, G, G_T, p) \leftarrow \text{MapGen}(1^\ell)$ , and  $g^\alpha, g^{\alpha^2}, \dots, g^{\alpha^q}$  where  $\alpha \stackrel{R}{\leftarrow} \mathbb{Z}_p$  and  $q$  is bounded by a polynomial in  $\ell$ , there exists no probabilistic polynomial-time algorithm  $\mathcal{A}$  that can compute  $(s, e(g, g)^{1/(\alpha+s)})$  where  $s \in \mathbb{Z}_p$  with a non-negligible probability in  $\ell$ . The probability is defined over the random choices of the parameters and random coins used by  $\mathcal{A}$ .

**Decisional Linear assumption (DL) [42].** For given  $(e, g, G, G_T, p) \leftarrow \text{MapGen}(1^\ell)$ , and  $(f, h, g^{r_1}, f^{r_2}, Q)$  where  $f, h, Q \stackrel{R}{\leftarrow} G$  and  $r_1, r_2 \stackrel{R}{\leftarrow} \mathbb{Z}_p$ , there exists no probabilistic polynomial-time algorithm  $\mathcal{A}$  that can determine  $Q \stackrel{?}{=} h^{r_1+r_2}$  with a non-negligible advantage, where “advantage” is defined as  $|\Pr[\mathcal{A}(g, f, h, g^{r_1}, f^{r_2}, Q) = 1] - \Pr[\mathcal{A}(g, f, h, g^{r_1}, f^{r_2}, h^{r_1+r_2}) = 1]|$ , and the probability is defined over the random choices of the parameters and random coins used by  $\mathcal{A}$ .

**Unforgeable Digital Signature.** Let  $\text{Sig} = (\text{sigKeyGen}, \text{sigSign}, \text{sigVerify})$  be a secure signature scheme,

Notation	Description
$D_a, C_a$ $D_b, C_b$	Alice's data set and its encryption form Bob's data set and its encryption form
Setup, KeyGen, Enc, Dec, AuGen, SetOp, Verify	algorithms of the VDSI
pm, sk, pk, si, au, rslt, proof	parameters of the VDSI
acKeyGen, acGen, acProve, acVerify	algorithms of the multi-accumulator
acSk, acPk, acDig, acRslt, acWit	parameters of the multi-accumulator
sigKeyGen, sigSign, sigVerify	algorithms of signature scheme Sig
sigSk, sigPk, $\sigma$	parameters of the signature scheme

TABLE I  
NOTATIONS FOR ALGORITHMS AND PARAMETERS IN THE VDSI,  
multi-accumulator AND THE SIGNATURE SCHEME Sig.

where sigKeyGen generates a pair of public and private keys, sigSign generates a signature for a message, and sigVerify determines if a message matches a signature. Any signature scheme satisfying the standard definition of unforgeability under adaptive chosen-message attacks [43] is sufficient for the purpose of this paper.

### III. VDSI MODEL AND DEFINITION

#### A. System Model

Figure 1 illustrates the *system model* of VDSI (verifiable delegated set intersection operations on outsourced encrypted data). The system has four entities: a trusted third party, a cloud, and two cloud users (i.e., data owners) referred to as Alice and Bob. The trusted third party is responsible for initializing system public parameters used by the cloud and cloud users. Alice and Bob can be either individuals or organizations that outsource their private data sets, denoted by  $D_a$  and  $D_b$ , to the cloud in encrypted form, denoted by  $C_a$  and  $C_b$ , respectively. Alice and Bob want to compute the intersection set  $D_a \cap D_b$ , by delegating the set intersection operation to the cloud but *without* giving the cloud the capability to decrypt  $C_a$  and  $C_b$ .

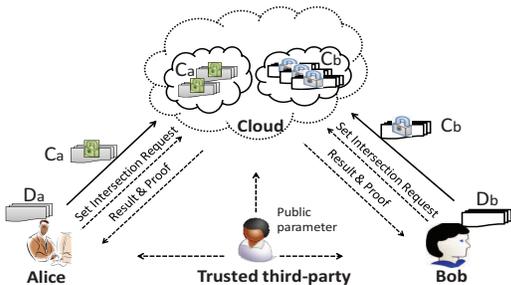


Fig. 1. VDSI system model: data owners Alice and Bob encrypt their data sets (denoted by  $D_a$  and  $D_b$ ), using their respective public keys, outsource to the cloud the resulting ciphertexts (denoted by  $C_a$  and  $C_b$ ), and delegate the computation of  $D_a \cap D_b$  to the cloud but without giving it the capability to decrypt  $C_a$  and  $C_b$ .

**Remark.** Note that data users have no prior knowledge about whom they will deal with for set intersection and therefore cannot share any common secret before outsourcing their data set, which rules out the many solutions with shared

secret [22]–[24]. In addition, the above system model can be easily extended to accommodate the following more general scenarios. First, rather than letting Alice and Bob outsource their encrypted data to the same cloud, they can outsource their encrypted data to two different clouds (dubbed *storage clouds*). Second, rather than letting (one of) the storage cloud(s) conduct the delegated computation on  $C_a$  and  $C_b$ , another cloud (dubbed *computing cloud*) or any other third party can be used for this purpose. The extension is trivial and omitted.

#### B. Threat Model and Basic Idea of Defense

We assume that the cloud users (i.e., data owners) are honest-but-curious, meaning that they act according to the protocols and use their real data sets as inputs to the protocols, but are curious about each other's private data. However, the cloud is possibly malicious. This means that the cloud can attempt to breach the secrecy of the data outsourced to the cloud, manipulate the integrity of the outsourced data, and deviate from the protocols arbitrarily. The cloud may be controlled by an attacker, who also has control over all the communication channels. This means that denial-of-service attack is inevitable and should be addressed orthogonally by another layer of defense. We will use “the attacker” and “the cloud” interchangeably.

The basic idea for defending against the possibly malicious cloud is to ask the cloud to generate a proof, which shows that it has faithfully executed the delegated set intersection operation. By “examining” the result and proof returned by the cloud, Alice and/or Bob can verify whether or not the cloud has faithfully executed the delegated set intersection operations on  $C_a$  and  $C_b$  or not.

#### C. VDSI Function Definition

In order to simplify the description of both the definition and the concrete scheme that will be presented later, we assume that there is an authenticated user-to-cloud private communication channel, which is used by a cloud user to send some secret information to the cloud (e.g., the secret information that allows the cloud to conduct the delegated set intersection operation). This assumption does not impose any significant restriction because in the case the cloud is controlled by the attacker, the secret information is given to the attacker anyway. In practice, the channel can be readily realized by encrypting the secret information under the cloud's public key.

Without loss of generality, denote by Alice's plaintext data set  $D_a = \{d_{a,0}, \dots, d_{a,n}\}$  and Bob's plaintext data set  $D_b = \{d_{b,0}, \dots, d_{b,m}\}$ . Alice (Bob) outsources her (his) encrypted version of  $D_a$  ( $D_b$ ), denoted by  $C_a$  ( $C_b$ ), to the cloud. Alice and Bob want to compute  $D_a \cap D_b$  by delegating the computation to the cloud, but without giving the cloud capability to decrypt  $C_a$  and  $C_b$ .

*Definition 1:* A VDSI scheme has seven algorithms:

- $\text{pm} \leftarrow \text{Setup}(1^\ell)$ : Given security parameter  $\ell$ , the trusted third party runs this algorithm to bootstrap the public parameters pm.

- $(pk_a, sk_a) \leftarrow \text{KeyGen}(pm)$ : Alice runs this randomized algorithm to generate a pair of public and private keys  $(pk_a, sk_a)$ , where  $pk_a$  is made public and  $sk_a$  is kept secret by Alice. Similarly, we denote Bob's pair of public and private keys by  $(pk_b, sk_b)$ .
- $(C_a, si_a) \leftarrow \text{Enc}(pk_a, D_a)$ : Alice runs this encryption algorithm to encrypt her data set  $D_a$  to ciphertext  $C_a$ , which is outsourced to the cloud, and some secret information  $si_a$ , which is kept secret by Alice. Bob can generate  $(C_b, si_b)$  similarly.
- $\{D', \perp\} \leftarrow \text{Dec}(sk_a, rslt_a)$ : Alice runs this decryption algorithm to decrypt ciphertext  $rslt_a$ , which is the output of the delegated set intersection operation conducted by the cloud on ciphertexts  $C_a$  and  $C_b$ , to obtain the intersection set  $D' = D_a \cap D_b$ . In the case the decryption fails, the algorithm outputs  $\perp$  instead. Note this decryption algorithm also can be used to decrypt  $C_a$ .
- $au_a \leftarrow \text{AuGen}(sk_a, si_a, pk_b)$ : In order to allow the cloud to conduct the set intersection operation on the outsourced ciphertexts  $C_a$  and  $C_b$ , Alice runs this algorithm to generate some auxiliary information  $au_a$ , which is sent to the cloud through the authenticated user-to-cloud private communication channel (see justification above), where  $si_a$  is the Alice's secret information generated by  $\text{Enc}(pk_a, D_a)$ . Similarly, Bob can generate and send  $au_b$  to the cloud, where  $au_b \leftarrow \text{AuGen}(sk_b, si_b, pk_a)$ .
- $\{(rslt_a, proof_a), (rslt_b, proof_b)\} \leftarrow \text{SetOp}(C_a, au_a, C_b, au_b)$ : This is the delegated set intersection operation run by the cloud. Depending on the application, the cloud may return  $(rslt_a, proof_a)$  and  $(rslt_b, proof_b)$  respectively to Alice and Bob, or return  $(rslt_a, proof_a)$  to Alice or  $(rslt_b, proof_b)$  to Bob who requested the delegated set intersection operation, where  $proof_a$  and  $proof_b$  are proofs that can show that the cloud has faithfully executed the SetOp protocol.
- $\{0, 1\} \leftarrow \text{Verify}(sk_a, si_a, rslt_a, proof_a)$ : Alice runs this algorithm to verify whether  $rslt_a$  is faithfully generated by the cloud according to the SetOp protocol. If so (with output 1), Alice calls the  $\text{Dec}(sk_a, rslt_a)$  algorithm to decrypt  $rslt_a$ ; otherwise (with output 0), the cloud is cheating.

We say a VDSI scheme is correct if the following holds:

$$\Pr \left[ \begin{array}{l} pm \leftarrow \text{Setup}(1^\ell), \\ (pk_a, sk_a) \leftarrow \text{KeyGen}(pm), \\ (pk_b, sk_b) \leftarrow \text{KeyGen}(pm), \\ \forall D_a, D_b, (C_a, si_a) \leftarrow \text{Enc}(pk_a, D_a), \\ (C_b, si_b) \leftarrow \text{Enc}(pk_b, D_b), \\ au_a \leftarrow \text{AuGen}(sk_a, si_a, pk_b), \\ au_b \leftarrow \text{AuGen}(sk_b, si_b, pk_a), \\ \{(rslt_a, proof_a), (rslt_b, proof_b)\} \leftarrow \\ \quad \text{SetOp}(C_a, au_a, C_b, au_b) : \\ 1 \leftarrow \text{Verify}(sk_a, si_a, rslt_a, proof_a), \\ 1 \leftarrow \text{Verify}(sk_b, si_b, rslt_b, proof_b), \\ D_a \cap D_b = \text{Dec}(sk_a, rslt_a) = \text{Dec}(sk_b, rslt_b) \end{array} \right] = 1.$$

#### D. VDSI Security Definition

Informally, VDSI aims to achieve the following security properties against the afore-discussed threat model. We consider three security properties: *outsourced data secrecy*, *function output secrecy* and *verifiability*, which are formally defined below. Let  $\epsilon$  be a negligible function in security parameter  $\ell$ . We consider a probabilistic polynomial-time (in  $\ell$ ) adversary  $\mathcal{A}$  controlling the cloud.

**Outsourced Data Secrecy:** Similar to security against chosen-plaintext attack, this property means that the attacker  $\mathcal{A}$  cannot breach secrecy of the outsourced data, unless that  $\mathcal{A}$  is provided with the auxiliary information.

*Definition 2:* (outsourced data secrecy) A VDSI scheme achieves *outsourced data secrecy* if the following holds:

$$\Pr \left[ \begin{array}{l} pm \leftarrow \text{Setup}(1^\ell), \\ (pk, sk) \leftarrow \text{KeyGen}(pm), \\ (D_0, D_1) \leftarrow \mathcal{A}^{\text{Enc}}(pk), s.t. |D_0| = |D_1|, \\ \lambda \xleftarrow{R} \{0, 1\}, (C_\lambda, si_\lambda) \leftarrow \text{Enc}(pk, D_\lambda), \\ \lambda' \leftarrow \mathcal{A}^{\text{Enc}}(pk, C_\lambda, D_0, D_1) : \\ \lambda = \lambda' \end{array} \right] - \frac{1}{2} \leq \epsilon$$

This property is necessary but not sufficient because it only assures the secrecy of outsourced data when  $\mathcal{A}$  is not given the delegated set intersection operation capability. The following property, function output secrecy, is used to capture the secrecy of outsourced data after  $\mathcal{A}$  is granted the capability (i.e.  $\mathcal{A}$  is given the auxiliary information  $au$ ).

**Function Output Secrecy:** This property means that  $\mathcal{A}$  cannot breach secrecy of the resulting intersection set  $D_a \cap D_b$ . Ideally, given a target ciphertext and the auxiliary information,  $\mathcal{A}$  cannot learn the plaintext with a non-negligible probability.

*Definition 3:* (function output secrecy) A VDSI scheme achieves *function output secrecy* if

$$\Pr \left[ \begin{array}{l} pm \leftarrow \text{Setup}(1^\ell), \\ (pk_a, sk_a) \leftarrow \text{KeyGen}(pm), \\ (pk_b, sk_b) \leftarrow \text{KeyGen}(pm), \\ \forall D_a, D_b, (C_a, si_a) \leftarrow \text{Enc}(pk_a, D_a), \\ (C_b, si_b) \leftarrow \text{Enc}(pk_b, D_b), \\ \forall cph \in (C_a \cup C_b), \\ au_a \leftarrow \text{AuGen}(sk_a, si_a, pk_b), \\ au_b \leftarrow \text{AuGen}(sk_b, si_b, pk_a), \\ \{d_1, \dots, d_q\} \leftarrow \mathcal{A}^{\text{Enc, SetOp, Verify}} \\ (pk_a, au_a, C_a, pk_b, au_b, C_b, cph) : \\ \exists i \in [1, q], d_i = d \end{array} \right] \leq f(\ell, q, |\mathcal{M}|),$$

where  $q$  is the maximum number of guessing against  $cph$ ,  $d$  is the plaintext with respect to  $cph$ , and  $\mathcal{M}$  is the plaintext domain.

**Remark.** Ideally, we want  $f(\ell, q, |\mathcal{M}|)$  to be a negligible function in  $\ell$  as well. Unfortunately, we are only able to construct a scheme that achieves  $f(\ell, q, |\mathcal{M}|) = \frac{q}{|\mathcal{M}|} + \epsilon$ , which is a non-negligible function in  $\ell$  because  $|\mathcal{M}|$  would not be exponentially in  $\ell$ . The intuition behind  $\frac{q}{|\mathcal{M}|}$  is that  $\mathcal{A}$  can launch a plaintext-guessing attack (in a way similar to the *online dictionary attack* against passwords), which is specific to our scheme that will be presented later. Designing a VDSI

scheme that achieves negligible  $f(\ell, q, |\mathcal{M}|)$  in  $\ell$  is left as an open problem for future research. Nevertheless, our definition is general enough to accommodate that scenario.

**Verifiability:** This property means that any  $\mathcal{A}$  not faithfully executing the SetOp protocol is bound to be caught.

*Definition 4:* (verifiability) A VDSI scheme is *verifiable* if

$$\Pr \left[ \begin{array}{l} \text{pm} \leftarrow \text{Setup}(1^\ell), \\ (\text{pk}_a, \text{sk}_a) \leftarrow \text{KeyGen}(\text{pm}), \\ (\text{pk}_b, \text{sk}_b) \leftarrow \text{KeyGen}(\text{pm}), \\ (D_a, D_b) \leftarrow \mathcal{A}^{\text{Enc, AuGen, SetOp, Verify}}(\text{pk}_a, \text{pk}_b), \\ (C_a, \text{si}_a) \leftarrow \text{Enc}(\text{pk}_a, D_a), \\ (C_b, \text{si}_b) \leftarrow \text{Enc}(\text{pk}_b, D_b), \\ \text{au}_a \leftarrow \text{AuGen}(\text{sk}_a, \text{si}_a, \text{pk}_b), \\ \text{au}_b \leftarrow \text{AuGen}(\text{sk}_b, \text{si}_b, \text{pk}_a), \\ \{(\text{rslt}_a, \text{proof}_a), (\text{rslt}_b, \text{proof}_b)\} \leftarrow \mathcal{A}^{\text{Enc, SetOp, Verify}} \\ (\text{pk}_a, \text{au}_a, D_a, C_a, \text{si}_a, \text{pk}_b, \text{au}_b, D_b, C_b, \text{si}_b) : \\ 1 \leftarrow \text{Verify}(\text{sk}_a, \text{si}_a, \text{rslt}_a, \text{proof}_a) \wedge \\ 1 \leftarrow \text{Verify}(\text{sk}_b, \text{si}_b, \text{rslt}_b, \text{proof}_b) \wedge \\ (\text{Dec}(\text{sk}_a, \text{rslt}_a) \neq \text{Dec}(\text{sk}_b, \text{rslt}_b)) \vee \\ \text{Dec}(\text{sk}_a, \text{rslt}_a) \neq (D_a \cap D_b) \end{array} \right]$$

#### IV. BUILDING-BLOCK: multi-accumulator

A cryptography accumulator is a primitive for a *verifier* to examine the membership of an element with respect to a (static or dynamic) data set. The examination is based on some public data and membership proof provided by a *prover*. In a single-accumulator scheme, each membership proof allows a verifier to examine the membership of a *single* element with respect to a data set. The idea of single-accumulator has been studied extensively (see, e.g., [44]–[46]). In this paper, we introduce the idea of multi-accumulator by which, each membership proof allows a verifier to examine the membership of *multiple* elements with respect to a data set. In the context of set intersection operations, multi-accumulator allows Alice (Bob) to verify, via a single examination, that  $D_a \cap D_b \subseteq D_a$  ( $\subseteq D_b$ ).

##### A. Function and Security Definitions

Suppose Alice has a data set  $\text{acD}_a$  and outsources it to the cloud (as the prover). Bob (as the verifier) has a dataset  $\text{acD}_b$  and queries the cloud for  $\text{acD}_a \cap \text{acD}_b$ .

*Definition 5:* A multi-accumulator scheme has the following algorithms:

- $(\text{acSk}, \text{acPk}) \leftarrow \text{acKeyGen}(1^\ell)$ : The trusted third party runs this algorithm to generate a pair of public and private key  $(\text{acPk}, \text{acSk})$ .
- $\text{acDig}_a \leftarrow \text{acGen}(\text{acPk}, \text{acD}_a)$ : Alice runs this algorithm to generate a digest  $\text{acDig}$  for  $\text{acD}_a$ , which is outsourced to the cloud. Similarly, Bob can generate  $\text{acDig}_b$  with respect to  $\text{acD}_b$ .
- $(\text{acRslt}, \text{acWit}) \leftarrow \text{acProve}(\text{acPk}, \text{acD}_b, \text{acD}_a)$ : Given the data set  $\text{acD}_b$  from Bob, the cloud runs this algorithm to generate  $\text{acRslt} = (\text{acD}_b \cap \text{acD}_a)$  with an accompanying witness  $\text{acWit}$  for this fact.
- $\{0, 1\} \leftarrow \text{acVerify}(\text{acPk}, \text{acDig}_b, \text{acRslt}, \text{acWit}, \text{acDig}_a)$ : Bob runs this algorithm to examine if  $\text{acRslt} = \text{acD}_b \cap \text{acD}_a$ , where  $\text{acDig}_b$  is the digest with respect to

$\text{acD}_b$  and  $\text{acDig}_a$  is the digest with respect to  $\text{acD}_a$ . If so, output 1; otherwise, output 0.

A multi-accumulator scheme is correct if

$$\Pr \left[ \begin{array}{l} \forall \text{acD}_a, \text{acD}_b \\ (\text{acSk}, \text{acPk}) \leftarrow \text{acKeyGen}(1^\ell), \\ \text{acDig}_b \leftarrow \text{acGen}(\text{acPk}, \text{acD}_b), \\ \text{acDig}_a \leftarrow \text{acGen}(\text{acPk}, \text{acD}_a), \\ (\text{acRslt}, \text{acWit}) \leftarrow \text{acProve}(\text{acPk}, \text{acD}_b, \text{acD}_a) : \\ 1 \leftarrow \text{acVerify}(\text{acPk}, \text{acDig}_b, \text{acRslt}, \text{acWit}, \text{acDig}_a) \end{array} \right] = 1.$$

A multi-accumulator scheme is secure if a malicious probabilistic polynomial-time prover  $\mathcal{A}$  can cheat the honest verifier without being caught. Let  $\ell$  be a security parameter and  $\epsilon$  be  $\leq \epsilon$  a negligible function in  $\ell$ . Formally, we have:

*Definition 6:* A multi-accumulator scheme is secure if

$$\Pr \left[ \begin{array}{l} (\text{acPk}, \text{acSk}) \leftarrow \text{acKeyGen}(1^\ell), \\ \text{acD}_a \leftarrow \mathcal{A}^{\text{acProve, acVerify}}(\text{acPk}), \\ \text{acDig}_a \leftarrow \text{acGen}(\text{acSk}, \text{acD}_a), \\ (\text{acD}_b, \text{acRslt}, \text{acWit}) \\ \leftarrow \mathcal{A}^{\text{acProve, acVerify}}(\text{acPk}, \text{acD}_a) : \\ \text{acDig}_b \leftarrow \text{acGen}(\text{acPk}, \text{acD}_b), \\ 1 \leftarrow \text{acVerify}(\text{acPk}, \text{acDig}_b, \text{acRslt}, \text{acWit}, \text{acDig}_a), \\ \text{acRslt} \neq \text{acD}_b \cap \text{acD}_a \end{array} \right] \leq \epsilon.$$

##### B. Construction based on Bilinear Map

A multi-accumulator scheme can be based on a single-accumulator scheme that supports both membership and non-membership proofs, as follows: the cloud generates a witness for each element of  $\text{acD}_b$  showing the element is a member or non-member of  $\text{acD}_a$  and simply puts them together as the witness for  $\text{acRslt} = \text{acD}_b \cap \text{acD}_a$ . However, this straightforward approach is costly because both the computational and communication complexities are linear to  $|\text{acD}_b|$ .

We present a multi-accumulator scheme, where the size of the witness is constant (i.e., independent of  $|\text{acD}_b|$ ). The proposed multi-accumulator scheme is extended from the single-accumulator scheme due to [44], [46], while adapting the basic idea underlying [16] as follows: Suppose Alice's data set is  $\text{acD}_a = \{d_{a,1}, \dots, d_{a,n}\}$ , Bob's data set is  $\text{acD}_b = \{d_{b,1}, \dots, d_{b,m}\}$ , and  $\text{acRslt} = \text{acD}_a \cap \text{acD}_b$ . We can encode  $\text{acD}_a$  via polynomial  $R(x) = \prod_{t \in \text{acD}_a} (x+t)$ , encode  $\text{acD}_b$  via polynomial  $W(x) = \prod_{t \in \text{acD}_b} (x+t)$ , encode the intersection set  $\text{acRslt}$  via polynomial  $T(x) = \prod_{t \in \text{acRslt}} (x+t)$ , and encode the subset  $\text{acD}_b - \text{acRslt}$  via polynomial  $Q(x) = \prod_{t \in (\text{acD}_b - \text{acRslt})} (x+t)$ . These polynomials satisfy the following: (i)  $T(x)Q(x) = W(x)$ , (ii)  $T(x)$  is a divisor of  $R(x)$ , and (iii)  $Q(x)$  is co-prime to  $R(x)$ . For the special case  $\text{acRslt} = \emptyset$ , the three conditions also hold since  $T(x) = 1$ ,  $Q(x) = W(x) = \prod_{t \in \text{acD}_b} (x+t)$  and  $R(x) = \prod_{t \in \text{acD}_a} (x+t)$ . Therefore, based on this idea, the multi-accumulator scheme allows the cloud to show the correctness of the intersection set, which can be either empty or non-empty. It can be constructed as follows:

- $\text{acKeyGen}(1^\ell)$ : Let  $(e, g, G, G_T, p) \leftarrow \text{MapGen}(1^\ell)$ , set  $\alpha \stackrel{R}{\leftarrow} \mathbb{Z}_p$  and  $\text{acPk} = (g^\alpha, g^{\alpha^2}, \dots, g^{\alpha^q})$ ,  $\text{acSk} = (\alpha)$ , where  $q$  is bounded by a polynomial in security parameter  $\ell$ .

- $\text{acGen}(\text{acPk}, \text{acD}_a)$ : Given Alice's data set  $\text{acD}_a = \{d_{a,1}, \dots, d_{a,n}\} \in \mathbb{Z}_p^n$  where  $n \leq q$ , compute its digest as

$$\text{acDig}_a = g^{\prod_{i=1}^n (d_{a,i} + \alpha)}.$$

- $\text{acProve}(\text{acPk}, \text{acD}_b, \text{acD}_a)$ : Given Bob's data set  $\text{acD}_b = (d_{b,1}, \dots, d_{b,m}) \in \mathbb{Z}_p^m$  where  $m \leq q$ , compute  $\text{acRslt} = \text{acD}_b \cap \text{acD}_a$ , and generate a witness as follows:

- Let  $T'(x) = \prod_{t \in (\text{acD}_a - \text{acRslt})} (x + t)$  and compute  $g^{T'(\alpha)}$  by substituting  $x$  with  $\alpha$ .
- Let  $Q(x) = \prod_{t \in (\text{acD}_b - \text{acRslt})} (x + t)$  and  $R(x) = \prod_{t \in \text{acD}_a} (x + t)$ , and find two polynomials  $Q'(x), R'(x)$  such that  $Q(x)Q'(x) + R(x)R'(x) = 1 \pmod p$  by taking advantage of  $\gcd(Q(x), R(x)) = 1$ . Compute  $(g^{Q(\alpha)}, g^{Q'(\alpha)}, g^{R(\alpha)}, g^{R'(\alpha)})$  by substituting  $x$  with  $\alpha$ .

Set  $\text{acRslt} = \text{acD}_b \cap \text{acD}_a$  and  $\text{acWit} = (g^{Q(\alpha)}, g^{Q'(\alpha)}, g^{R(\alpha)}, g^{R'(\alpha)})$ .

- $\text{acVerify}(\text{acPk}, \text{acD}_b, \text{acRslt}, \text{acWit}, \text{acDig}_a)$ : Given  $\text{acWit}$  and  $\text{acRslt}$  from the prover, the verifier proceeds as follows:

- 1) If  $\text{acRslt} \neq \emptyset$ , compute  $g^{T(\alpha)}$  according to  $T(x) = \prod_{t \in \text{acRslt}} (x + t)$ . Otherwise, let  $T(x) = 1$  and  $g^{T(\alpha)} = g$ .
- 2) If  $e(g^{Q(\alpha)}, g^{T(\alpha)}) \neq e(\text{acDig}_b, g)$ , return 0; otherwise, proceed to next step.
- 3) If  $e(g^{T(\alpha)}, g^{T'(\alpha)}) \neq e(\text{acDig}_a, g)$ , return 0; otherwise, proceed to next step.
- 4) If  $e(g^{Q(\alpha)}, g^{Q'(\alpha)})e(\text{acDig}_a, g^{R'(\alpha)}) \neq e(g, g)$ , return 0; otherwise, return 1.

Correctness of the multi-accumulator scheme can be verified easily. We describe its asymptotic efficiency in Table II. It is worth noting that (i) the witness generated by algorithm  $\text{acProve}$  only consists of four group elements, meaning that the complexity is independent of  $k = |\text{acD}_b \cap \text{acD}_a|$ , and (ii) the computational complexity of algorithm  $\text{acVerify}$  is linear to  $k = |\text{acD}_b \cap \text{acD}_a|$ .

TABLE II

ASYMPTOTICAL EFFICIENCY OF THE multi-accumulator SCHEME, WHERE  $\text{Exp}$  DENOTES THE EXPONENTIATION OPERATION,  $\text{Pairing}$  DENOTES THE PAIRING OPERATION,  $n = |\text{acD}_a|$ ,  $m = |\text{acD}_b|$  AND  $k = |\text{acD}_a \cap \text{acD}_b|$ .

	acGen	acProve	acVerify
Computation	$n\text{Exp}$	$(n + m)\text{Exp}$	$k\text{Exp} + 7\text{Pairing}$
Output Size	$ G $	$4 G $	N/A

The security of the multi-accumulator can be assured by the following theorem, the proof of which is shown in [47].

*Theorem 1:* Assume that the  $q$ -SDH assumption holds, the multi-accumulator scheme is secure with respect to Definition 6.

## V. THE VDSI SCHEME

**Basic Ideas.** In order to attain a VDSI scheme, we need to resolve two issues: (i) How can we enable the cloud to compare the equality of two ciphertexts that are encrypted under two different public keys  $\text{pk}_a$  and  $\text{pk}_b$ , respectively?

(ii) How can we enable the cloud to generate a proof for showing that it has faithfully executed the SetOp protocol, ideally without using zero-knowledge proof for the sake of better efficiency?

To resolve the above (i), we adopt the idea of proxy re-encryption as follows: Alice can generate a re-key and send it to the cloud, which can use the re-key to transform ciphertext  $C_a$  (encrypted under Alice's public key  $\text{pk}_a$ ) into an intermediate form, say  $T_a$ . Similarly, the cloud can transform ciphertext  $C_b$  (encrypted under Bob's public key  $\text{pk}_b$ ) into the same kind intermediate form, denoted by  $T_b$ . Then, the cloud can "compare"  $T_a$  and  $T_b$  to determine whether they correspond to the same plaintext or not. More specifically, a data item  $d_{a,i} \in D_a$  is encrypted using  $\text{pk}_a = (g^{\beta_a}, g^{\gamma_a})$  as  $(g^{r_2}, g^{\gamma r_1}, d_{a,i} g^{\beta(r_1+r_2)})$ , where  $r_1, r_2 \xleftarrow{R} \mathbb{Z}_p$ . Alice can give the re-key  $\text{rk}_a = g^{\beta_a/\gamma_a}$ , rather than her private key  $\text{sk}_a = (\beta_a, \gamma_a)$  to the cloud, which now can transform the ciphertext into

$$\frac{e(d_{a,i} g^{\beta_a(r_1+r_2)}, g)}{e(g^{\gamma_a r_1}, g^{\beta_a/\gamma_a}) e(g^{r_2}, g^{\beta_a})} = e(d_{a,i}, g).$$

Similarly, for data item  $d_{b,i} \in D_b$ , the cloud can transform the corresponding ciphertext into  $e(d_{b,i}, g)$ . If  $d_{a,i} = d_{b,i}$ , then  $e(d_{a,i}, g) = e(d_{b,i}, g)$ . While this method is sufficient to allow the cloud to determine whether the two ciphertexts correspond to the same plaintext or not, it does not achieve the desired semantic security because the cloud can launch the *plaintext guess* attack against elements  $d_{a,i}$  and  $d_{b,i}$ . We have tried without success to eliminate this attack while preserving the other properties (especially the *verifiability*). We therefore leave it as an open problem.

To resolve the above issue (ii), we observe that the cloud, as illustrated above, can generate  $e(d_{a,i}, g)$  for each  $d_{a,i} \in D_a$  and  $e(d_{b,i}, g)$  for each  $d_{b,i} \in D_b$ . As a result, the cloud can use the multi-accumulator scheme to generate a proof as follows: For Bob, let  $e(d_{a,i}, g)$ 's as  $\text{acD}_a$  and  $e(d_{b,i}, g)$ 's as  $\text{acD}_b$ , the cloud applies  $\text{acProve}$  to generate witness  $\text{acWit}_b$  for showing  $\text{acRslt} = \text{acD}_a \cap \text{acD}_b$  is the correct intersection set with respect to  $e(d_{a,i}, g)$ 's and  $e(d_{b,i}, g)$ 's. Given witness  $\text{acWit}_b$ , digest  $\text{acDig}_b$  of  $e(d_{b,i}, g)$ 's and digest  $\text{acDig}_a$  of  $e(d_{a,i}, g)$ , Bob can verify the correctness of  $\text{acRslt}$ , which can be computed from the returned intersection set of  $C_a$  and  $C_b$ . Similarly, the cloud can generate witness  $\text{acWit}_a$  for Alice, who can then conduct the same kind of verification. That is, by using the multi-accumulator scheme in section IV, the cloud users can verify the correctness of the intersection set, which contains zero or more common elements.

### A. The Scheme

The scheme is a modular construction based on (i) a secure multi-accumulator scheme  $\text{Ac} = (\text{acKeyGen}, \text{acGen}, \text{acProve}, \text{acVerify})$  such as the one described in Section IV, and (ii) a secure digital signature scheme  $\text{Sig} = (\text{sigKeyGen}, \text{sigSign}, \text{sigVerify})$ . The digital signature scheme is used to authenticate the encryption form of the accumulator digest, which assures that the cloud cannot manipulate it without being detected. Specifically, the scheme is described as follows:

**Setup**( $1^\ell$ ): Given security parameter  $\ell$ , the trusted third party runs  $(e, g, G, G_T, p) \leftarrow \text{MapGen}(1^\ell)$ . Let  $H : G_T \rightarrow \mathbb{Z}_p$  be a collision-resistant hash function. The trusted third party also runs  $(\text{acPk}, \text{acSk}) \leftarrow \text{KeyGen}(1^\ell)$  and sets the public parameter as

$$\text{pm} = (\text{acPk}, e, p, g, G, G_T).$$

**KeyGen**( $\text{pm}$ ): Alice runs  $(\text{sigPk}_a, \text{sigSk}_a) \leftarrow \text{sigKeyGen}(1^\ell)$ , selects  $\beta_a, \gamma_a \xleftarrow{R} \mathbb{Z}_p$ , and sets

$$\text{sk}_a = (\beta_a, \gamma_a, \text{sigSk}_a), \quad \text{pk}_a = (g^{\beta_a}, g^{\gamma_a}, \text{sigPk}_a).$$

Similarly, Bob generates  $\text{sk}_b = (\beta_b, \gamma_b, \text{sigSk}_b)$  and  $\text{pk}_b = (g^{\beta_b}, g^{\gamma_b}, \text{sigPk}_b)$ .

**Enc**( $\text{pk}_a, D_a$ ): Alice, with  $D_a = (d_{a,1}, \dots, d_{a,n})$  where  $d_{a,i} \in G$  for  $1 \leq i \leq n$ , executes as follows:

- Select  $d_{a,0} \xleftarrow{R} G$  (for a security purpose that will be elaborated later).
- For  $0 \leq i \leq n$ , select  $r_{i1}, r_{i2} \xleftarrow{R} \mathbb{Z}_p$  and compute

$$\text{cph}_{a,i} = (g^{r_{i2}}, g^{\gamma_a r_{i1}}, d_{a,i} g^{\beta_a (r_{i1} + r_{i2})}).$$

- For  $0 \leq i \leq n$ , let  $T_i = H(e(d_{a,i}, g))$  and compute  $\text{acDig}_a \leftarrow \text{acGen}(\text{acPk}, \{T_0, \dots, T_n\})$ .
- Set  $C_a = \{\text{cph}_{a,0}, \dots, \text{cph}_{a,n}\}$  and  $\text{si}_a = \text{acDig}_a$ .

Similarly, Bob, with  $D_b = (d_{b,1}, \dots, d_{b,m})$  where  $d_{b,i} \in G$  for  $1 \leq i \leq m$ , can obtain  $C_b = \{\text{cph}_{b,0}, \dots, \text{cph}_{b,m}\}$  and  $\text{si}_b = \text{acDig}_b$ .

**Dec**( $\text{sk}_a, \text{rslt}_a$ ): Given the cloud-generated ciphertext intersection set  $\text{rslt}_a = \{\text{cph}_{a,j}, \dots, \text{cph}_{a,k}\}$  where  $1 \leq j, k \leq n$ , Alice decrypts ciphertexts  $\text{cph}_{a,i}$  for  $j \leq i \leq k$  as follows:

$$d_{a,i} = d_{a,i} g^{\beta_a (r_{i1} + r_{i2})} / (g^{r_{i2}})^{\beta_a} (g^{\gamma_a r_{i1}})^{\beta_a / \gamma_a}.$$

The decryption of  $\text{rslt}_a$  is  $D_a \cap D_b = \{d_{a,j}, \dots, d_{a,k}\}$ . Note that this algorithm can also be used to decrypt  $C_a$  without involving any delegated set operations. In this case, the integrity of  $C_a$  can be easily assured by  $\text{acDig}_a$  since the plaintexts of  $C_a$  should be accumulated to  $\text{acDig}_a$ .

Similarly, Bob can decrypt the cloud-generated ciphertext intersection set  $\text{rslt}_b = \{\text{cph}_{b,j}, \dots, \text{cph}_{b,k}\}$  where  $1 \leq j, k \leq m$  to obtain  $D_a \cap D_b$ .

**AuGen**( $\text{sk}_a, \text{si}_a, \text{pk}_b$ ): Given private key  $\text{sk}_a$ , Alice generates re-key  $\text{rk}_a = (g^{\beta_a / \gamma_a})$ . Alice encrypts the secret information  $\text{si}_a$  using Bob's public key  $\text{pk}_b$  to obtain ciphertext  $\text{cph}_B = (g^{r_2}, g^{\gamma_b r_1}, \text{acDig}_a g^{\beta_b (r_1 + r_2)})$ , where  $r_1, r_2 \xleftarrow{R} \mathbb{Z}_p$ . Then, Alice runs  $\sigma_a \leftarrow \text{sigSign}(\text{sigSk}_a, \text{cph}_B)$  to obtain a signature  $\sigma_a$  on message  $\text{cph}_B$ . Finally, Alice sets  $\text{au}_a = (\text{rk}_a, \text{cph}_B, \sigma_a)$ .

Similarly, Bob can generate  $\text{au}_b = (\text{rk}_b, \text{cph}_A, \sigma_b)$ .

**SetOp**( $C_a, \text{au}_a, C_b, \text{au}_b$ ): Given  $C_a = \{\text{cph}_{a,0}, \dots, \text{cph}_{a,n}\}$ ,  $C_b = \{\text{cph}_{b,0}, \dots, \text{cph}_{b,m}\}$ ,  $\text{au}_a = (\text{rk}_a = g^{\beta_a / \gamma_a}, \text{cph}_B, \sigma_a)$ , and  $\text{au}_b = (\text{rk}_b = g^{\beta_b / \gamma_b}, \text{cph}_A, \sigma_b)$ , the cloud executes as follows:

- Transform ciphertexts  $\text{cph}_{a,i}$  for  $0 \leq i \leq n$  into

$$T_{a,i} = \frac{e(d_{a,i} g^{\beta_a (r_{i1} + r_{i2})}, g)}{e(g^{\gamma_a r_{i1}}, g^{\beta_a / \gamma_a}) e(g^{r_{i2}}, g^{\beta_a})} = e(d_{a,i}, g),$$

and compute  $T_a = \{H(T_{a,0}), \dots, H(T_{a,n})\}$ .

- Transform ciphertexts  $\text{cph}_{b,i}$  for  $1 \leq i \leq m$  into

$$T_{b,i} = \frac{e(d_{b,i} g^{\beta_b (r_{i1} + r_{i2})}, g)}{e(g^{\gamma_b r_{i1}}, g^{\beta_b / \gamma_b}) e(g^{r_{i2}}, g^{\beta_b})} = e(d_{b,i}, g)$$

and compute  $T_b = \{H(T_{b,0}), \dots, H(T_{b,m})\}$ .

- Generate the intersection set  $\text{rslt}_a$  and a proof with respect to  $C_a$  as follows: Run  $(\text{acRslt}, \text{acWit}_a) \leftarrow \text{acProve}(\text{acPk}, T_a, T_b)$  and set

$$\begin{aligned} \text{rslt}_a &= \{\text{cph}_{a,i} | H(A_{a,i}) \in \text{acRslt}\}, \\ \text{proof}_a &= (\text{acWit}_a, \text{cph}_A, \sigma_b). \end{aligned}$$

- Generate the intersection set  $\text{rslt}_b$  and a proof with respect to  $C_b$  as follows: Run  $(\text{acRslt}, \text{acWit}_b) \leftarrow \text{acProve}(\text{acPk}, T_b, T_a)$  and set

$$\begin{aligned} \text{rslt}_b &= \{\text{cph}_{b,i} | H(A_{b,i}) \in \text{acRslt}\}, \\ \text{proof}_b &= (\text{acWit}_b, \text{cph}_B, \sigma_a). \end{aligned}$$

**Verify**( $\text{sk}_a, \text{si}_a, \text{rslt}_a, \text{proof}_a$ ): Given  $\text{rslt}_a$  and  $\text{proof}_a$ , Alice verifies that the cloud faithfully executed the SetOp protocol as follows:

- Verify the integrity of  $\text{cph}_A$  by running  $\text{sigVerify}(\text{sigPk}_b, \text{cph}_A, \sigma_b)$ . If it outputs 0, then return 0; otherwise, proceed to next step.
- Decrypt  $\text{cph}_A$  using private key  $\text{sk}_a$  according to

$$\text{acDig}_b = \text{acDig}_b g^{\beta_b (r_1 + r_2)} / (g^{r_2})^{\beta_b} (g^{\gamma_a r_1})^{\beta_b / \gamma_a}.$$

- If  $\text{rslt}_a$  is not empty, decrypt  $\text{rslt}_a$  to obtain the plaintexts and compute  $Y_a = \{e(d_{a,i}, g) | \text{cph}_{a,i} \in \text{rslt}_a\}$ . Otherwise, let  $Y_a = \emptyset$ .
- Run  $\text{acVerify}(\text{acPk}, \text{acDig}_a, Y_a, \text{acWit}_a, \text{acDig}_b)$ . If it outputs 0, then return 0; otherwise, return 1.

If the algorithm returns 1,  $\text{Dec}(\text{sk}_a, \text{rslt}_a)$  is called to obtain  $D_a \cap D_b$ .

Similarly, Bob can run the same algorithm to verify that the cloud does not cheat.

**Remark: why using  $d_{a,0}$  and  $d_{b,0}$ ?** Since Alice needs to know the accumulator digest  $\text{acDig}_b$  for the sake of verifying the correctness of  $\text{rslt}_a$ , we need to assure that Alice cannot use  $\text{acDig}_b$  to infer useful information about  $D_b$ . This is achieved by "blending" the accumulator digest with the randomness, namely the hash value of the randomly selected  $d_{b,0}$ . This eliminates the usage of zero-knowledge proofs [48], [49], while assuring no useful information is leaked.

**Remark.** Our VDSI scheme only offers coarse-grained access control in the following sense. Suppose Alice and Bob allow the cloud to conduct the delegated set intersection operation on  $C_a$  and  $C_b$ , and Alice and Carlos allow the cloud to conduct the delegated set intersection operation on  $C_a$  and  $C_c$ . Then, the cloud is able to conduct the set intersection operation on  $C_b$  and  $C_c$  without the authorization from Bob and Carlos. It is a future work to enforce fine-grained access control in VDSI.

## B. Security Analysis

Correctness of the VDSI scheme can be examined easily. The security properties can be assured by the following theorems, the proofs of which can be found in [47].

*Theorem 2:* Under the DL assumption, the scheme achieves *outsourced data secrecy* (Definition 2).

*Theorem 3:* Given that the bilinear map  $e$  is one-way, the VDSI scheme achieves the function output secrecy property (Definition 3).

*Theorem 4:* Assume that Sig is an unforgeable signature scheme, Ac is a secure multi-accumulator scheme and  $H$  is a collision resistance hash function, the VDSI scheme achieves the verifiability property (Definition 4).

## VI. PERFORMANCE EVALUATION

### A. Asymptotic Complexity

Table III summarizes the communication and computational overhead incurred by the VDSI solution, which is grouped into two phases: (i) data outsourcing phase, during which the cloud users outsource their encrypted private data sets to the cloud (i.e. Enc); (ii) set operation phase, during which the cloud users attain the intersection set (i.e. AuGen, SetOp and Verify). We compare the overhead with its counterpart that is incurred by the straightforward solution, namely that Alice and Bob download their outsourced encrypted data from the cloud, decrypt their data, and then run a PSI protocol to jointly compute the intersection set. From the perspective of cloud users, we observe that the VDSI scheme outperforms the straightforward solution in both communication and computational complexities. Assume that the data sets have been stored in the cloud, the VDSI scheme only incurs  $O(k)$  computational complexity to obtain the intersection set (including the cost for verification) and  $O(k)$  communication complexity for returning the intersection set to the user, where  $k$  is the size of intersection set. This means that the VDSI scheme is optimal (up to a constant factor). In contrast, the straightforward solution incurs  $O(m + n)$  in computational and communication overhead where  $m$  and  $n$  are the sizes of the two data sets. The advantage of VDSI scheme is most substantial when  $k \ll m$  or  $k \ll n$ .

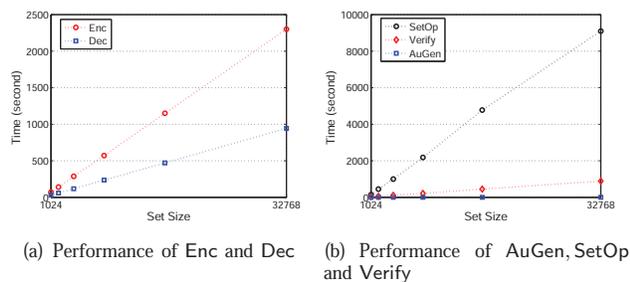


Fig. 2. Performance of VDSI, where Alice and Bob outsource their data sets of the same size (i.e.,  $m = n$ ), algorithms Enc, Dec, AuGen and Verify run on the CLIENT MACHINE, and algorithm SetOp runs on the SERVER MACHINE.

## B. Performance Evaluation

**Implementation** We implemented the VDSI scheme in JAVA based on the Java Pairing Based Cryptography library (jPBC) [50]. In our implementation, we instantiated the bilinear map with Type A pairing ( $\ell = 512$ ), which offers a level of security that is comparable to 1024-bit DLOG [50]. We instantiated the signature scheme with the DSA signature scheme provided by JDK1.6. We varied the set size ( $m$  and  $n$ ) from  $2^{10}$  to  $2^{15}$ . The algorithms run by the cloud users (i.e., Enc, Dec, AuGen and Verify) were executed on a CLIENT MACHINE with Linux OS, 2.93GHz Intel Core Duo CPU (E7500), and 2GB RAM. The algorithm run by the cloud (i.e., SetOp) was executed on a SERVER MACHINE with Linux OS, 4 processors of 2.40GHz Intel Xeon CPU, and 8GB RAM.

**Evaluation and result** In our experiments we set the same size for data sets owned by the two cloud users, i.e.,  $m = n$ , and set the size of intersection set  $k = n/2$ . For algorithms Enc, Dec, AuGen and Verify, we evaluated each algorithm's execution time for both cloud users and treat their average execution time as the real execution time. Figure 2(a) plots the execution time of Enc and Dec that are run by the cloud users. We observe that the execution times for both algorithms are almost linear to the size of data sets. We also can see that the execution of Enc is more expensive than that of Dec. However, Enc is executed only once when the cloud user outsources data sets. Figure 2(b) shows the execution time of SetOp (run by the cloud) and the execution time of AuGen and Verify (run by the cloud users). We observe that the execution time of AuGen and Verify is much more smaller than that of the algorithm SetOp. This suggests that cloud users should leverage the cloud's computation resources by delegating set intersection operations.

**Performance Comparison** In order to understand the benefit and limitation of the VDSI solution, we compare it with the straightforward solution, where data sets are encrypted by the algorithm Enc and Dec of the VDSI, and the private set operations between two data users are performed by the protocol (Java version) in [19], which is the most efficient PSI protocol in the literature.

We ran the experiments on the same SERVER MACHINE with Linux OS, 4 processors of 2.40GHz Intel Xeon CPU, and 8GB RAM, with the data sets each consisting of 32768 elements. We vary the size of the intersection set as 25%, 50% and 75% of the size of the data set respectively, and compare the communication and computation overhead in the set operation phase (we did not compare the cost of the data outsourcing phase because they are the same for both solutions), which is shown in Figure 3. From Figure 3(a) we observe that the computation overhead in the VDSI solution decreases when the size of the intersection set decreases. However, the computation overhead in the straightforward solution remains the same regardless the size of the intersection set. We also can see in Figure 3(b) that the communication overhead for the data users in the VDSI solution is linear to the size of intersection set, and is much less than that of the straightforward solution. This advantage can become

TABLE III

ASYMPTOTIC PERFORMANCE COMPARISON FOR THE VDSI SCHEME AND THE STRAIGHTFORWARD SOLUTION. WE ASSUME THAT THE STRAIGHTFORWARD SOLUTION ADOPTING THE ENCRYPTION AND DECRYPTION ALGORITHMS OF THE VDSI SCHEME. HERE  $n$  IS THE SIZE OF ALICE'S DATA SET,  $m$  IS THE SIZE OF BOB'S DATA SET,  $k$  IS THE SIZE OF SET INTERSECTION, AND  $\text{Comp}(\text{PSI})$  AND  $\text{Comm}(\text{PSI})$  DENOTES THE RESPECTIVE COMPUTATION AND COMMUNICATION COMPLEXITY OF THE PRIVATE SET INTERSECTION PROTOCOL. NOTE THAT  $\text{Comp}(\text{PSI})$  AND  $\text{Comm}(\text{PSI})$  ARE BOTH LINEAR TO THE SIZE OF DATA SETS ( $m + n$ ) FOR THE STATE-OF-THE-ART SOLUTION [19].

Phase		VDSI solution			Straightforward solution		
		Alice	Bob	Cloud	Alice	Bob	Cloud
Data outsourcing	Computation	$O(n)$	$O(m)$	N/A	$O(n)$	$O(m)$	N/A
	Communication	$O(n)$	$O(m)$	$O(n + m)$	$O(n)$	$O(m)$	$O(n + m)$
Set operation	Computation	$O(k)$	$O(k)$	$O(m + n)$	$O(n) + \text{Comp}(\text{PSI})$	$O(m) + \text{Comp}(\text{PSI})$	N/A
	Communication	$O(k)$	$O(k)$	$O(k)$	$O(n) + \text{Comm}(\text{PSI})$	$O(m) + \text{Comm}(\text{PSI})$	$O(n + m)$

more substantial when the size of the intersection size is far less than the size of data set owned by the data users. We note that while the straightforward solution can use other efficient encryption schemes (e.g., symmetric encryption) to encrypt/decrypt data sets to achieve higher efficiency, it cannot reduce the communication cost that is linear to the size of data set. Therefore, our VDSI solution is extremely suitable for computing intersection set whose size is far less than that of the data sets.

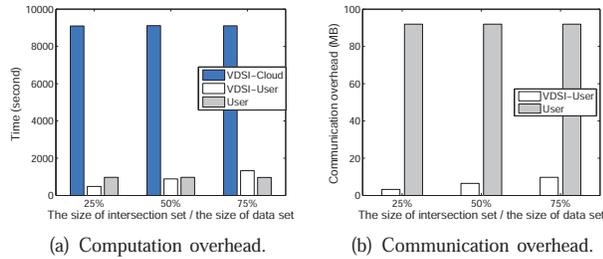


Fig. 3. Performance comparison between our VDSI solution and the straightforward solution, where each data user outsourced the data set of 32768 elements. We vary the size of the intersection set with 25%, 50% and 75% of the size of the data set respectively. VDSI-User and VDSI-Cloud denote the costs spent by the cloud and each data user in the VDSI solution and User represents the cost spent by each data user in the straightforward solution.

### C. Improvement with Parallelization

In our proposed scheme, the execution of Enc, Dec, SetOp and Verify can be implemented more efficiently with parallelization, because operations related to elements of data sets are independent. In practice, we implemented the algorithms by using multiple threads to compute independent operations (e.g. encrypting elements of the data sets, decrypting ciphertexts, and transforming ciphertexts into a value of  $G_T$ ). In the parallelization version, we created 4 threads and ran the algorithms Enc, Dec, SetOp and Verify on the SERVER MACHINE with Linux OS, 4 processors of 2.40GHz Intel Xeon CPU, and 8GB RAM. To understand the efficiency gain of parallelization, we also ran the algorithms without parallelization on the same SERVER MACHINE. Figure 4 shows the performance comparison, which indicates that the algorithms using parallelization are about 2 times faster than their counterparts that do not use parallelization. This means that our scheme can leverage the multi-core architecture, and

that our scheme is suitable for delegating set intersection over outsourced large data sets.

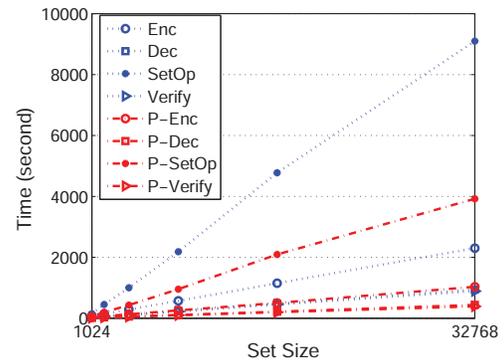


Fig. 4. Performance comparison for algorithms Enc, Dec, AuGen and Verify executed on SERVER MACHINE. The algorithms with prefix "P-" were implemented with parallelization, and the algorithms without prefix were implemented without parallelization.

## VII. CONCLUSION

We have introduced the novel notion of VDSI, which allows two users to outsource to the cloud their encrypted data sets as well as the set intersection operation on ciphertexts. This is achieved without giving the cloud the capability to decrypt the encrypted data, while enabling the users to hold the misbehaving cloud accountable.

Our study brings interesting and challenging open problems for future research. In addition to the ones mentioned in the paper (e.g., incorporating fine-grained access control, if possible), we need to design the same kinds of solutions for other set operations.

## ACKNOWLEDGMENT

This work was partly supported by NSF under Grant No. 1111925.

## REFERENCES

- [1] C. Gentry, "Computing arbitrary functions of encrypted data," *Commun. ACM*, vol. 53, no. 3, pp. 97–105, Mar. 2010. [Online]. Available: <http://doi.acm.org/10.1145/1666420.1666444>
- [2] Z. Brakerski, C. Gentry, and V. Vaikuntanathan, "(leveled) fully homomorphic encryption without bootstrapping," in *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, ser. ITCS '12. New York, NY, USA: ACM, 2012, pp. 309–325. [Online]. Available: <http://doi.acm.org/10.1145/2090236.2090262>

- [3] C. Gentry, S. Halevi, and N. P. Smart, "Fully homomorphic encryption with polylog overhead," in *Proceedings of the 31st Annual International Conference on Theory and Applications of Cryptographic Techniques*, ser. EUROCRYPT'12. Berlin, Heidelberg: Springer-Verlag, 2012, pp. 465–482. [Online]. Available: [http://dx.doi.org/10.1007/978-3-642-29011-4\\_28](http://dx.doi.org/10.1007/978-3-642-29011-4_28)
- [4] K. Lauter, M. Naehrig, and V. Vaikuntanathan, "Can homomorphic encryption be practical?" Cryptology ePrint Archive, Report 2011/405, 2011, <http://eprint.iacr.org/>.
- [5] P. Baldi, R. Baronio, E. D. Cristofaro, P. Gasti, and G. Tsudik, "Countering gattaca: efficient and secure testing of fully-sequenced human genomes," in *ACM Conference on Computer and Communications Security*, 2011, pp. 691–702.
- [6] E. D. Cristofaro and G. Tsudik, "Practical private set intersection protocols with linear computational and bandwidth complexity," *IACR Cryptology ePrint Archive*, vol. 2009, p. 491, 2009.
- [7] G. Mezzour, A. Perrig, V. Gligor, and P. Papadimitratos, "Privacy-preserving relationship path discovery in social networks," in *Cryptology and Network Security*, ser. Lecture Notes in Computer Science, vol. 5888. Springer Berlin Heidelberg, 2009, pp. 189–208.
- [8] G. Ateniese, E. D. Cristofaro, and G. Tsudik, "(if) size matters: Size-hiding private set intersection," in *Public Key Cryptography*, 2011, pp. 156–173.
- [9] M. J. Freedman, K. Nissim, and B. Pinkas, "Efficient private matching and set intersection," in *EUROCRYPT*, 2004, pp. 1–19.
- [10] L. Kissner and D. X. Song, "Privacy-preserving set operations," in *CRYPTO*, 2005, pp. 241–257.
- [11] J. Camenisch and G. M. Zaverucha, "Private intersection of certified sets," in *Financial Cryptography*, 2009, pp. 108–127.
- [12] B. Pinkas, T. Schneider, and M. Zohner, "Faster private set intersection based on OT extension," in *Proceedings of the 23rd USENIX Security Symposium, San Diego, CA, USA, August 20-22, 2014.*, 2014, pp. 797–812.
- [13] A. E. Kosba, D. Papadopoulos, C. Papamanthou, M. F. Sayed, E. Shi, and N. Triandopoulos, "TRUESET: faster verifiable set computations," in *Proceedings of the 23rd USENIX Security Symposium, San Diego, CA, USA, August 20-22, 2014.*, 2014, pp. 765–780.
- [14] C. Hazay and Y. Lindell, "Efficient protocols for set intersection and pattern matching with security against malicious and covert adversaries," *J. Cryptology*, vol. 23, no. 3, pp. 422–456, 2010.
- [15] D. Dachman-Soled, T. Malkin, M. R. 0001, and M. Yung, "Efficient robust private set intersection," in *ACNS*, 2009, pp. 125–142.
- [16] C. Papamanthou, R. Tamassia, and N. Triandopoulos, "Optimal verification of operations on dynamic sets," Cryptology ePrint Archive, Report 2010/455, 2010, <http://eprint.iacr.org/>.
- [17] C. Hazay and K. Nissim, "Efficient set operations in the presence of malicious adversaries," in *Public Key Cryptography*, 2010, pp. 312–331.
- [18] S. Jarecki and X. Liu, "Efficient oblivious pseudorandom function with applications to adaptive ot and secure computation of set intersection," in *TCC*, 2009, pp. 577–594.
- [19] C. Dong, L. Chen, and Z. Wen, "When private set intersection meets big data: An efficient and scalable protocol," Cryptology ePrint Archive, Report 2013/515, 2013, <http://eprint.iacr.org/>.
- [20] Y. Huang, D. Evans, and J. Katz, "Private set intersection: Are garbled circuits better than custom protocols?" *19th Network and Distributed Security Symposium*, 2012.
- [21] E. D. Cristofaro and G. Tsudik, "Practical private set intersection protocols with linear complexity," in *Financial Cryptography*, 2010, pp. 143–159.
- [22] F. Kerschbaum, "Collusion-resistant outsourcing of private set intersection," in *SAC*, 2012, pp. 1451–1456.
- [23] —, "Outsourced private set intersection using homomorphic encryption," in *ASIACCS*, 2012, pp. 85–86.
- [24] S. Kamara, P. Mohassel, M. Raykova, and S. Sadeghian, "Server-aided private set intersection: Scaling to million element sets," 2013, <http://research.microsoft.com/pubs/194141/sapsi.pdf>.
- [25] R. Canetti, O. Paneth, D. Papadopoulos, and N. Triandopoulos, "Verifiable set operations over outsourced databases," Cryptology ePrint Archive, Report 2013/724, 2013, <http://eprint.iacr.org/>.
- [26] G. Yang, C. H. Tan, Q. Huang, and D. S. Wong, "Probabilistic public key encryption with equality test," in *CT-RSA* 2010, pp. 119–131.
- [27] S. Canard, G. Fuchsbaauer, A. Gouget, and F. Laguillaumie, "Plaintext-checkable encryption," in *CT-RSA*, 2012, pp. 332–348.
- [28] B. Wang, M. Li, S. Chow, and H. Li, "Computing encrypted cloud data efficiently under multiple keys," in *Communications and Network Security (CNS), 2013 IEEE Conference on*, Oct 2013, pp. 504–513.
- [29] M. Abadi, D. Boneh, I. Mironov, A. Raghunathan, and G. Segev, "Message-locked encryption for lock-dependent messages," in *CRYPTO (I)*, 2013, pp. 374–391.
- [30] Q. Tang, "Towards public key encryption scheme supporting equality test with fine-grained authorization," in *ACISP*, 2011, pp. 389–406.
- [31] B. Parno, M. Raykova, and V. Vaikuntanathan, "How to delegate and verify in public: Verifiable computation from attribute-based encryption," in *TCC*, 2012, pp. 422–439.
- [32] S. Goldwasser, Y. T. Kalai, and G. N. Rothblum, "Delegating computation: interactive proofs for muggles," in *STOC*, 2008, pp. 113–122.
- [33] R. Gennaro, C. Gentry, and B. Parno, "Non-interactive verifiable computing: Outsourcing computation to untrusted workers," in *CRYPTO*, 2010, pp. 465–482.
- [34] K.-M. Chung, Y. T. Kalai, and S. P. Vadhan, "Improved delegation of computation using fully homomorphic encryption," in *CRYPTO*, 2010, pp. 483–501.
- [35] C. Papamanthou, E. Shi, and R. Tamassia, "Signatures of correct computation," in *TCC*, 2013, pp. 222–242.
- [36] S. Benabbas, R. Gennaro, and Y. Vahlis, "Verifiable delegation of computation over large datasets," in *CRYPTO*, 2011, pp. 111–131.
- [37] D. Fiore and R. Gennaro, "Publicly verifiable delegation of large polynomials and matrix computations, with applications," in *Proceedings of the 2012 ACM conference on Computer and communications security*, ser. CCS '12. New York, NY, USA: ACM, 2012, pp. 501–512. [Online]. Available: <http://doi.acm.org/10.1145/2382196.2382250>
- [38] R. Gennaro, C. Gentry, B. Parno, and M. Raykova, "Quadratic span programs and succinct nizks without pcps," in *EUROCRYPT*, 2013, pp. 626–645.
- [39] Q. Chai and G. Gong, "Verifiable symmetric searchable encryption for semi-honest-but-curious cloud servers," in *ICC*, 2012, pp. 917–922.
- [40] B. Thompson, S. Haber, W. G. Horne, T. Sander, and D. Yao, "Privacy-preserving computation and verification of aggregate queries on outsourced databases," in *Privacy Enhancing Technologies*, 2009, pp. 185–201.
- [41] L. Zhang, X.-Y. Li, Y. Liu, and T. Jung, "Verifiable private multi-party computation: Ranging and ranking," in *INFOCOM*, 2013, pp. 605–609.
- [42] D. Boneh, X. Boyen, and H. Shacham, "Short group signatures," in *In proceedings of CRYPTO 04, LNCS series*. Springer-Verlag, 2004, pp. 41–55.
- [43] S. Goldwasser, S. Micali, and R. L. Rivest, "A digital signature scheme secure against adaptive chosen-message attacks," *SIAM J. Comput.*, vol. 17, no. 2, pp. 281–308, Apr. 1988. [Online]. Available: <http://dx.doi.org/10.1137/0217017>
- [44] L. Nguyen, "Accumulators from bilinear pairings and applications," in *Proceedings of the 2005 international conference on Topics in Cryptology*, ser. CT-RSA'05. Berlin, Heidelberg: Springer-Verlag, 2005, pp. 275–292. [Online]. Available: [http://dx.doi.org/10.1007/978-3-540-30574-3\\_19](http://dx.doi.org/10.1007/978-3-540-30574-3_19)
- [45] J. Camenisch and A. Lysyanskaya, "Dynamic accumulators and application to efficient revocation of anonymous credentials," in *Proceedings of the 22nd Annual International Cryptology Conference on Advances in Cryptology*, ser. CRYPTO '02. London, UK: Springer-Verlag, 2002, pp. 61–76. [Online]. Available: <http://dl.acm.org/citation.cfm?id=646767.704437>
- [46] I. Damgård and N. Triandopoulos, "Supporting non-membership proofs with bilinear-map accumulators," *IACR Cryptology ePrint Archive*, vol. 2008, p. 538, 2008.
- [47] Q. Zheng and S. Xu, "Verifiable delegated set intersection operations on outsourced encrypted data," Cryptology ePrint Archive, Report 2014/178, 2014, <http://eprint.iacr.org/>.
- [48] U. Feige, A. Fiat, and A. Shamir, "Zero-knowledge proofs of identity," *J. Cryptol.*, vol. 1, no. 2, pp. 77–94, Aug. 1988. [Online]. Available: <http://dx.doi.org/10.1007/BF02351717>
- [49] T. Acar, S. S. M. Chow, and L. Nguyen, "Accumulators and u-prove revocation," in *Financial Cryptography*, 2013, pp. 189–196.
- [50] "The java pairing based cryptography library." <http://gas.dia.unisa.it/projects/jpbcl/>.