

Linear Private Set Union from Multi-Query Reverse Private Membership Test

Cong Zhang^{1,2}, Yu Chen^{3,4,5} (✉), Weiran Liu⁶, Min Zhang^{3,4,5} and Dongdai Lin^{1,2}

¹State Key Laboratory of Information Security, Institute of Information Engineering, Chinese Academy of Sciences, Beijing 100093, China

²School of Cyber Security, University of Chinese Academy of Sciences, Beijing 100049, China
{zhangcong, ddlin}@iie.ac.cn

³School of Cyber Science and Technology, Shandong University, Qingdao 266237, China

⁴State Key Laboratory of Cryptology, P.O. Box 5159, Beijing 100878, China

⁵Key Laboratory of Cryptologic Technology and Information Security, Ministry of Education, Shandong University, Qingdao 266237, China
yuchen.prc@gmail.com, zm_min@mail.sdu.edu.cn

⁶Alibaba Group

weiran.lwr@alibaba-inc.com

Abstract

Private set union (PSU) protocol enables two parties, each holding a set, to compute the union of their sets without revealing anything else to either party. So far, there are two known approaches for constructing PSU protocols. The first mainly depends on additively homomorphic encryption (AHE), which is generally inefficient since it needs to perform a non-constant number of homomorphic computations on each item. The second is mainly based on oblivious transfer and symmetric-key operations, which is recently proposed by Kolesnikov et al. (ASIACRYPT 2019). It features good practical performance, which is several orders of magnitude faster than the first one. However, neither of these two approaches is optimal in the sense that their computation and communication complexity are not both $O(n)$, where n is the size of the set. Therefore, the problem of constructing the optimal PSU protocol remains open.

In this work, we resolve this open problem by proposing a generic framework of PSU from oblivious transfer and a newly introduced protocol called multi-query reverse private membership test (mq-RPMT). We present two generic constructions of mq-RPMT. The first is based on symmetric-key encryption and general 2PC techniques. The second is based on re-randomizable public-key encryption. Both constructions lead to PSU with linear computation and communication complexity.

We implement our two PSU protocols and compare them with the state-of-the-art PSU. Experiments show that our PKE-based protocol has the lowest communication of all schemes, which is $3.7 - 14.8\times$ lower depending on set size. The running time of our PSU scheme is $1.2 - 12\times$ faster than that of state-of-the-art depending on network environments.

1 Introduction

Private set union (PSU) enables two parties, each holding a private set of elements, to compute the union of the two sets while revealing nothing more than the union

itself. PSU and its variants have numerous applications [7, 18, 24, 27, 29, 32, 33, 43]. An important PSU application is IP blacklist and vulnerability data aggregation [24, 43]. Consider two organizations (i.e. the maintainers of the IP blacklists) want to compute their IP blacklist joint list, which will help minimize vulnerabilities in their infrastructure. However, it is not secure to let the organizations simply exchange their blacklists because each individual IP blacklist is generated according to the detection strategy formulated by the maintainer and cannot be leaked. Note that a curious organization may infer the detection strategy of another organization from the IP address in the intersection. Therefore, it is important to hide the intersection, which is exactly the functionality of PSU.

Another killer application of PSU is to construct Private-ID protocol [8, 18]. The Private-ID protocol enables two parties, each holding a private set of items, to privately compute a set of random universal identifiers (UID) corresponding to the records in the union of their sets, where each party additionally learns which UIDs correspond to which items in its set but not if they belong to the intersection or not. The main use of Private ID is to realize data alignment, that is, both parties can sort their private data according to these universal identifiers. They can then proceed item-by-item, doing any desired private computation. Garimella et al. [18] gave a modular way to construct Private ID from Oblivious PRF (OPRF) and PSU. Their experiments showed that the bottleneck of their Private ID is the underlining PSU instantiations.

In addition, PSU applications also include information security risk assessment [33], joint graph computation [7], distributed network monitoring [29], building block for private DB supporting full join [32] etc.

Over the last decade, there has been a significant amount of work on private set operation, especially private set intersection (PSI) [11, 15, 31, 37, 38, 40]. We refer the reader to [41] for an overview of different PSI paradigms. State-of-the-art semi-honest PSI protocols in the two-party setting [11, 19, 31, 37, 44] all mainly rely on symmetric-key operations, except for a few base OT operations in OT extension

protocol [26, 30]. Let n denote the size of input set, the communication complexity of these OT-based PSI protocols has been improved from initial nonlinear $O(n \log n)$ [31, 39, 40] to linear complexity $O(n)$ [11, 14, 19, 20, 37, 44].

1.1 Motivation

In contrast to the affairs of PSI, the efficiency of the state-of-the-art PSU is less satisfactory. Roughly, there are two known approaches for constructing PSU protocols. The first is mainly based on public-key techniques. Existing constructions along this approach [16, 23, 29, 46] have to perform a non-constant number of additively homomorphic encryption (AHE) operations on each set element, rendering the overall protocols inefficient. The other is mainly based on symmetric-key techniques in combination with OT [18, 27, 32], which is several orders of magnitude faster than AHE-based constructions. However, neither of the two approaches is optimal in the sense that their computation and communication complexity are not both $O(n)$, where n is the size of the set. We note that [12] is the work closest to optimal bound, but its communication and computation complexity additionally depend on the statistical security parameter λ . This leaves the following open problem:

Can we construct PSU protocols with linear computation and communication complexity?

1.2 Our Contribution

In this paper, we answer this question affirmatively in the semi-honest setting. Our contribution can be summarized as follows:

1. We revisit the PSU protocol [32] (KRTW protocol for short hereafter) in depth. Roughly, KRTW protocol is built upon two building blocks, namely oblivious transfer (OT) and reverse private membership test (RPMT). We figure out the root causing KRTW protocol non-optimal is that RPMT has linear communication complexity and super-linear computation complexity, and it has to be carried out n times independently, where n is the size of sender’s private set.
2. To achieve linear complexity, we propose a new framework for constructing PSU protocols. The core building block is a newly introduced protocol called multi-query RPMT (mq-RPMT). We identify and overcome several technical difficulties for building optimal mq-RPMT, and give two realizations of mq-RPMT. Both the two concrete mq-RPMT protocols achieve linear communication and computation complexity.
3. We further abstract a new primitive called membership encryption (ME), which broadens the scope of the candidate encryption scheme, unifies our two constructions,

and halves the communication complexity of our SKE-based construction on receiver side.

4. Combining OT and the above mq-RPMT, we eventually obtain SKE-based and PKE-based PSU protocols with optimal complexity for the first time. Experiments show that our PKE-based protocol has the lowest communication of all schemes, which is $3.7 - 14.8\times$ lower depending on set size. The running time of our PSU scheme is $1.2 - 12\times$ faster than that of state-of-the-art depending on network environments. In addition to our scheme, we also use Silent OT [6, 48] to optimize the scheme of [18, 27], and provide different parameter selection of Ferret OT [48].

Figure 1 depicts the technical overview of our new PSU framework. We elaborate the details in the next subsection.

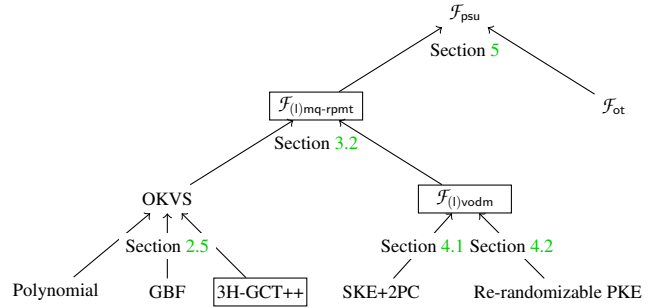


Figure 1: Technical overview of our new PSU framework. The new primitives and functionalities are marked with rectangles.

1.3 Overview of Our Techniques

We provide the high-level technical overview for our new framework of PSU protocol.

KRTW protocol revisit. Our starting point is the recent PSU protocol of Kolesnikov et al. [32]. The core of KRTW protocol is a subprotocol called reverse private membership test (RPMT), which can test whether a sender’s element y belongs to the receiver’s input set X , and let the receiver obtain the result. After that, both parties execute OT protocol to let the receiver obtain $\{y\} \cup X$. The computation cost of original RPMT [32] is $O(n \log^2 n)$ and the communication cost is $O(n)$. For the purpose of computing the set union, the parties need to execute RPMT n times independently, which results in $O(n^2)$ communication and $O(n^2 \log^2 n)$ computation. The complexity can be further reduced to $O(n \log n)$ and $O(n \log n \log \log n)$ separately via hash to bin technology, but it is still *super-linear*. The bottleneck of the KRTW protocol is exactly RPMT.

Zoom in on the original RPMT. The original RPMT protocol employs an oblivious PRF (OPRF) functionality $\mathcal{F}_{\text{opr}}^{\text{opr}}$ and a private equality test (PEQT) functionality $\mathcal{F}_{\text{peqt}}$. In

OPRF, the sender learns a random PRF key k and the receiver learns the PRF output $F_k(y_1), \dots, F_k(y_n)$ on its inputs $y_1, \dots, y_n \in Y$. In PEQT, the functionality receives two strings from the receiver and the sender respectively and tells the receiver whether the two strings are equal. Their RPMT uses an *indication string* s to indicate the membership of X .

More precisely, their RPMT protocol executes as follows with sender \mathcal{S} 's input y and receiver \mathcal{R} 's input $X = \{x_1, \dots, x_n\}$: \mathcal{S} and \mathcal{R} execute the OPRF protocol first. The receiver \mathcal{R} receives a PRF key k . The sender \mathcal{S} inputs y , and receives $q^* = F_k(y)$. Next, \mathcal{R} chooses a random indication string s . Then, \mathcal{R} computes and sends the interpolation polynomial P which passes through points $\{(x_i, s \oplus F_k(x_i))\}_{i \in [n]}$ to the sender. After receiving P , \mathcal{S} computes $s^* := q^* \oplus P(y)$. Now, \mathcal{S} and \mathcal{R} invoke the $\mathcal{F}_{\text{peqt}}$ -functionality with input s^* and s separately. Finally, \mathcal{R} receives output from $\mathcal{F}_{\text{peqt}}$.

If $y \in X$, i.e., there exists an x_i such that $y = x_i$, then we have $s^* = q^* \oplus P(y) = F_k(x_i) \oplus P(x_i) = s$. If $y \notin X$, then $q^* = F_k(y)$ is pseudorandom, which implies that $s^* = q^* \oplus P(y) \neq s$ with overwhelming probability.

To identify the root of the inefficiency of the original RPMT protocol, we first try to interpret it at an abstract level. Our first key observation is that the polynomial actually plays the role of oblivious key-value store (OKVS). Our second key observation is that the usage of OPRF is three-fold. Firstly, \mathcal{R} uses an OPRF to derive n pseudorandom one-time pads, then encrypts the same indication string into n ciphertexts under these one-time pads. Secondly, \mathcal{S} utilizes OPRF to decrypt a ciphertext obliviously. Finally, OPRF provides OKVS with randomness to ensure the correctness of the protocol.

Based on the above new interpretation, we are ready to describe our new mq-RPMT protocol in an incremental way over the original RPMT protocol.

Enhanced oblivious key-value store. One reason that accounts for the super-linear complexity of the original RPMT protocol is that the polynomial related operations are costly. More precisely, the complexity of polynomial interpolation is $O(n \log^2 n)$, and the amortized complexity of polynomial evaluation is $O(\log^2 n)$. According to our first observation, polynomial essentially plays the role of OKVS. This greatly increases the space of the concrete mapping schemes that can be used. A drop-in replacement of polynomial with more efficient OKVS candidates can reduce the computation complexity immediately. However, as we observed before, an additional randomness property should be satisfied now, since we do not use OPRF to provide randomness anymore. To achieve this goal, we enhance OKVS in two aspects: efficiency and security. (See Section 2.5 for the details.)

Oblivious decryption-then-matching. Another reason that accounts for the super-linear complexity is that the original RPMT protocol is one-time in nature. To see this, note that in the original RPMT protocol \mathcal{S} learns the purported indication string. This design lets \mathcal{S} learn more information than needed, and is exactly the reason that hinders multi-query. For exam-

ple, if there are two distinct elements belonging to \mathcal{R} 's set, then \mathcal{S} will obtain the same indication string. This will let \mathcal{S} know that the two elements belong to the intersection, which violates security.

Based on the above discussion, the rough idea of making RPMT support multi-query is to encode the *ciphertext* of indication string in OKVS instead of the indication string itself. In this way, \mathcal{S} will obtain some ciphertexts (i.e. the value of $\text{OKVS}(y)$), and \mathcal{R} has the corresponding key. We need to let \mathcal{R} decrypt these ciphertexts, and match the results with the indication string. A naive attempt is to have \mathcal{S} directly send the ciphertexts to \mathcal{R} , and in the sequel, \mathcal{R} tries to decrypt and match. However, this rough idea is problematic since it is insecure even against a semi-honest receiver. Consider \mathcal{R} records the correspondence between x_i and $\text{OKVS}(x_i)$. In this way, \mathcal{R} is able to learn \mathcal{S} 's private input y by simple look-up when $y \in X$, rather than merely the fact that $y \in X$. We overcome this difficulty in two steps. The first step is to re-factor the functionality of OPRF to encryption and oblivious decryption functionality. Let \mathcal{R} encrypt the indication string locally. Then \mathcal{R} computes the corresponding OKVS and sends it to \mathcal{S} . To ensure the overall protocol still constitutes an RPMT protocol, the second step is to merge the oblivious decryption functionality and PEQT into a new functionality, namely, vector oblivious decryption-then-matching (VODM) functionality. In this functionality, the sender inputs a vector of ciphertexts and the receiver inputs a key and a plaintext. The functionality decrypts these ciphertexts with the key and matches the results with the plaintext input by the receiver. If it matches, the receiver outputs 1, and outputs 0 otherwise.

Putting all the pieces together, we can build mq-RPMT protocol from OKVS, encryption, and VODM functionality in a modular way. (See Section 3 for the technical details).

Two generic constructions of mq-RPMT. Our first generic construction chooses probabilistic SKE as the encryption scheme, and resorts to general 2PC to implement the VODM functionality. See Section 4.1 for details. Our second generic construction chooses re-randomizable PKE as the encryption scheme and uses re-randomization technique to implement VODM functionality, without resorting to generic 2PC.

Our idea is to let \mathcal{S} re-randomize all the ciphertexts and then send the results to \mathcal{R} . In this way, \mathcal{R} fulfills the decryption-then-matching functionality in an oblivious manner for all $y_i \in X$. We note that this method will leak some information of $y \notin X$, however, as observed by KRTW, this leakage does not cause any harm to the PSU, since the PSU protocol releases that value anyway.

Looking ahead, one may doubt our PKE-based scheme is inefficient. We note that our PKE-based scheme can still be very efficient because we use PKE techniques in an entirely different way compared to prior PKE-based protocol [12, 16, 29]. We only need to perform the encryption, re-randomization, and decryption operations per item, while they need to carry

out many ciphertext homomorphism operations per item. See Section 4.2 for details.

Optimization with membership encryption. In the above framework, the underlying encryption schemes must be probabilistic to make the security proof go through. As a consequence, this incurs considerable overhead on communication costs due to ciphertext expansion. Observe that the VODM functionality reveals only one-bit information for every ciphertext. A second thought indicates that a full-fledged encryption scheme might be overkill for our construction of mq-RPMT protocol, and a new type of encryption scheme suffices. We propose the new encryption scheme as membership encryption (ME).

We sketch the definition of ME in the symmetric key setting as below. Let X be a string set. The encryption algorithm takes a key k and an element $x_i \in X$ as inputs, outputs a ciphertext c . The decryption algorithm takes a key k and a ciphertext c as inputs, outputs a bit to indicate if the encrypted element belongs to X . For the correctness, we require that for any $x_i \in X$ and any $c \leftarrow \text{Enc}(k, x_i)$, we have $\text{Dec}(k, c) = 1$. The security requirement is multi-element pseudorandomness, namely, $\{\text{Enc}(k, x_i)\}_{x_i \in X}$ are computationally indistinguishable to C^n , i.e. the uniform distribution over ciphertext space. The consistency requirement is that a random ciphertext decrypts to “0” with overwhelming probability.

Membership encryption distills the right functionality we need for an encryption scheme in mq-RPMT protocol. It not only encompasses the constructions from randomized SKE and PKE in a unified manner, but also admits new construction from deterministic SKE, which enjoys compact ciphertext. As we elaborate in Section 4.3, this new construction helps to halve the communication complexity on the receiver side.

1.4 Related Work

We survey existing PSU protocols with security against semi-honest adversaries. Hereafter, unless otherwise declared, we calculate the efficiency by assuming a balanced setting, namely the sets of both sender and receiver are of size n .

Kissner and Song [29] proposed the first PSU protocol based on polynomial representations and additively homomorphic encryption (AHE). The polynomial representation of a set is to represent a set by a polynomial f , in which each set item is the root of the polynomial. The main observation of them is that when the set of two parties is represented by polynomials f and g , the root of $f \cdot g$ is exactly the union items. The communication and computation complexity of the protocol are both quadratic to the set size n , and the efficiency is very low. Later, Frikken [16] found that it is enough to represent only receiver’s set in polynomial f . Then the receiver sends the AHE of f to the sender. The sender computes and sends back the ciphertexts of $(f(y), yf(y))$ for all $y \in Y$. In this way, the receiver could decrypted these ciphertexts and obtained the element outside of his set, since the intersection elements

decrypted to 0. Davidson and Cid [12] proposed a similar PSU protocol like Frikken, the main difference is that they use Bloom Filter (BF) instead of polynomial to represent the set. Both their protocols are expensive due to the frequent uses of AHE. Kolesnikov et al. [32] proposed the first PSU protocol mainly based on symmetric key techniques, which makes several orders of magnitude improvement of PSU. Recently, Garimella et al. [18] and Jia et al. [27] both use the oblivious switching network (OSN) subprotocol [35] to construct PSU, which further improve $2 - 4\times$ over [32]. However, all these symmetric key based PSU have the superlinear complexity.

Other PSU protocols focus on multi-party settings [5, 25, 29, 45], malicious settings [16, 23, 45] and computation with untrusted third party’s help [9, 10, 46]. All of the above constructions rely heavily on expensive AHE or zero-knowledge proof, which are out of the scope of our consideration.

Table 1 provides an asymptotic comparison of our design with the previous PSU works. We note that although the complexity of our SKE-based scheme is also related to t , where t is the number of AND gates in an SKE decryption circuit, we emphasize that t is a constant which is independent of n , that is, t remains the same no matter how n changes. In this sense, the complexity of our SKE-based scheme is strictly linear in n , though in practice t is larger than $\log n$. We leave the construction of a linear SKE-based PSU with a concrete complexity smaller than $\log n$ to future work.

Protocol	Communication	Computation
[29]	$O(\kappa^3 n^2)$	$O(n^2)$ pub
[16]	$O(\kappa n)$	$O(n \log \log n)$ pub
[12]	$O(\kappa \lambda n)$	$O(\lambda n)$ pub
[32]	$O(\kappa n \log n)$	$O(n \log n \log \log n)$ sym
[18]	$O(\kappa n \log n)$	$O(n \log n)$ sym
[27]	$O(\kappa n \log n)$	$O(n \log n)$ sym
Our SKE-based	$O((\kappa + t)n)$	$O(tn)$ sym
Our PKE-based	$O(\kappa n)$	$O(n)$ pub

Table 1: Asymptotic communication and computation costs of two-party PSU protocols in the semi-honest setting.

Pub: public-key operations; sym: symmetric cryptographic operations. n is the size of the parties’ input sets. κ and λ are computational and statistical security parameter respectively (typically $\kappa = 128$ and $\lambda=40$). t is the number of AND gates in an SKE decryption circuit. We ignore the pub-key cost of κ base OTs.

2 Preliminaries

Full version of this paper. Due to space constraints, we defer details like instantiation details, omitted proofs, omitted protocols, implementation details and supplementary experiments to the full version of this paper [50].

2.1 Notation

We denote the parties as receiver \mathcal{R} and sender \mathcal{S} , and their respective input sets as X and Y with $|X| = n_x$ and $|Y| = n_y$. In the balanced setting, we often just assume that $n = n_x = n_y$. We use κ and λ to denote the computational and statistical security parameters, respectively. We use $[n]$ to denote the set $\{1, 2, \dots, n\}$. For a bit string v we let v_i denote the i th bit. We use \mathbb{F}_{2^σ} to denote finite field composed of all σ -long bit strings. We say that a function f is negligible in κ if it vanishes faster than the inverse of any polynomial in κ , and write it as $f(\kappa) = \text{negl}(\kappa)$. We use the abbreviation PPT to denote probabilistic polynomial-time. By $a \stackrel{r}{\leftarrow} A$, we denote that a is randomly selected from the set A , $a \leftarrow A(x)$ denotes that a is the output of the randomized algorithm A on input x , and $a := b$ denotes that a is assigned by b .

2.2 Security Model

This work, similar to most protocols for private set operation, operates in the *semi-honest model*, where adversaries may try to learn as much information as possible from a given protocol execution but are not able to deviate from the protocol steps. This is in contrast to malicious adversaries which are able to deviate arbitrarily from the protocol. PSU protocols for the malicious setting exist, e.g., [5, 16, 23, 29, 45], but they are less efficient than protocols for the semi-honest setting.

Semi-honest security. We use the standard security definition for two-party computation [21] in this work.

Definition 2.1 Let $\text{view}_S^\Pi(X, Y)$ and $\text{view}_R^\Pi(X, Y)$ be the views of \mathcal{S} and \mathcal{R} in the protocol, and let $\text{output}(X, Y)$ be the output of both parties in protocol. A protocol Π is said to securely compute functionality f in the semi-honest model if for every PPT adversary \mathcal{A} there exists a PPT simulator Sim_S and Sim_R such that for all inputs X and Y ,

$$\{\text{view}_S^\Pi(X, Y), \text{output}(X, Y)\} \approx_c \{\text{Sim}_S(X, f(X, Y)), f(X, Y)\}$$

$$\{\text{view}_R^\Pi(X, Y), \text{output}(X, Y)\} \approx_c \{\text{Sim}_R(Y, f(X, Y)), f(X, Y)\}$$

2.3 Encryption Schemes

Our construction requires some encryption schemes. We use the standard definition of symmetric-key encryption (SKE) and re-randomizable public-key encryption (ReRand-PKE) schemes. For our purpose, we require a case-tailored security notion called *single-message multi-ciphertext pseudorandomness*. We give these definitions in the full version.

2.4 Oblivious Transfer

Oblivious transfer (OT) [42] is an important cryptographic primitive used in various multiparty computation protocols.

We define the functionality of 1-out-of-2 OT in Figure 2.

Parameters: Sender \mathcal{S} , Receiver \mathcal{R} , message length κ

Functionality:

- Wait for input $b \in \{0, 1\}$ from the receiver \mathcal{R} .
- Wait for input (x_0, x_1) from the sender \mathcal{S} .
- Give x_b to the receiver \mathcal{R} .

Figure 2: 1-out-of-2 Oblivious Transfer Functionality \mathcal{F}_{ot}

2.5 Oblivious Key-Value Stores

A key-value store [19, 38] is simply a data structure that maps a set of keys to corresponding values. The definition is as follows:

Definition 2.2 (Key-Value Store) A key-value store is parameterized by a set \mathcal{K} of keys, a set \mathcal{V} of values, and a set of function H , and consists of two algorithms:

- $\text{Encode}_H(\{(x_1, y_1), \dots, (x_n, y_n)\})$: on input key-value pairs $\{(x_i, y_i)\}_{i \in [n]} \subseteq \mathcal{K} \times \mathcal{V}$, outputs an object D (or, with statistically small probability, an error \perp).
- $\text{Decode}_H(D, x)$: on input D and a key x , outputs a value $y \in \mathcal{V}$.

Correctness. For all $A \subseteq \mathcal{K} \times \mathcal{V}$ with distinct keys:

$$(x, y) \in A \text{ and } \perp \neq D \leftarrow \text{Encode}_H(A) \implies \text{Decode}_H(D, x) = y$$

Obliviousness. For all distinct $\{x_1^0, \dots, x_n^0\}$ and all distinct $\{x_1^1, \dots, x_n^1\}$, if Encode_H does not output \perp for $\{x_1^0, \dots, x_n^0\}$ or $\{x_1^1, \dots, x_n^1\}$, then the distribution of $\{D | y_i \leftarrow \mathcal{V}, i \in [n], \text{Encode}_H((x_1^0, y_1), \dots, (x_n^0, y_n))\}$ is computationally indistinguishable to $\{D | y_i \leftarrow \mathcal{V}, i \in [n], \text{Encode}_H((x_1^1, y_1), \dots, (x_n^1, y_n))\}$.

A key-value store is an oblivious key-value store (OKVS) if it satisfies the obliviousness property.

Intuitively, obliviousness means that when value is randomly selected, the distribution of D is independent from key's set. In addition, our application requires OKVS to meet the *Randomness* property defined below to argue the correctness of our scheme.

Randomness. For any $A = \{(x_1, y_1), \dots, (x_n, y_n)\}$ and $x^* \notin \{x_1, \dots, x_n\}$, the output of $\text{Decode}_H(D, x^*)$ is statistically indistinguishable to that of uniform distribution over \mathcal{V} , where $D \leftarrow \text{Encode}_H(A)$.

The efficiency of an OKVS scheme can be measured by following three parameters:

- **Rate:** Let ratio n/m be the rate of key-value store, where m is the size of object D . Note that the optimal rate is 1.
- **Encoding complexity:** The computational cost of the Encode_H algorithm, as a function of the number n of key-value pairs.

- **Decoding complexity:** The computational cost of the Decode_H algorithm.

We investigated the existing schemes and found that the main candidates for OKVS are: Polynomial, Garbled Bloom Filter (GBF) [13] and Garbled Cuckoo Table (GCT) [19, 38, 44] etc. We give the general introduction and detailed comparisons of above OKVS in the full version.

Before instantiation, 3H-GCT recently proposed by Garimella et al. [19] could be a good candidate, which has linear encoding complexity $O(n)$ and a rate of 0.81. However, we find that the original 3H-GCT did not meet the *Randomness* we defined before because it was set to 0 in some positions of D . To solve this problem, a natural idea is to set random values in these positions like [44] does. We call this modified 3H-GCT as 3H-GCT++. We give the formal description of 3H-GCT++ in Figure 3 and we give the proof that our 3H-GCT++ satisfies obliviousness and randomness in the full version.

2.6 Private Set Union

PSU is a special case of secure two-party computation. The ideal functionality for PSU is given in Figure 4.

3 Multi-Query Reverse Private Membership Test

3.1 Definition

We propose mq-RPMT and give the formal definition of mq-RPMT functionality in Figure 5. For generality we set $|Y| = n_y$ and $|X| = n_x$ in our definition.

We define the vector oblivious decryption-then-matching $\mathcal{F}_{\text{vodm}}$ corresponding to encryption scheme \mathcal{E} in Figure 6, as a component of mq-RPMT.

3.2 Framework of Multi-Query RPMT

Now we describe our framework of mq-RPMT protocol. As we said in Section 1.3, the cryptographic primitives we use are a single-message multi-ciphertext pseudorandomness encryption scheme $\mathcal{E} = (\text{Setup}, \text{KeyGen}, \text{Enc}, \text{Dec})$, an OKVS scheme $(\text{Encode}_H, \text{Decode}_H)$ and the $\mathcal{F}_{\text{vodm}}$ functionality.

Let $Y = \{y_1, \dots, y_{n_y}\}$ and $X = \{x_1, \dots, x_{n_x}\}$ be the input of mq-RPMT sender \mathcal{S} and receiver \mathcal{R} . First, the receiver \mathcal{R} picks an indication string s^1 . Then \mathcal{R} chooses a random key k used in encryption scheme \mathcal{E} to encrypt s for n_x times, and obtains (s_1, \dots, s_{n_x}) . Next, \mathcal{R} computes an OKVS $D := \text{Encode}_H((x_1, s_1), \dots, (x_{n_x}, s_{n_x}))$ and sends D to \mathcal{S} . After receiving D , \mathcal{S} computes $s_i^* = \text{Decode}_H(D, y_i)$ for $i \in [n_y]$. Now \mathcal{S} and \mathcal{R} invoke the VODM functionality $\mathcal{F}_{\text{vodm}}$. \mathcal{S} acts

¹In fact, our indication string s could be any fixed value, e.g. $s = 0$, while s in KRTW must be selected randomly.

as sender with input $S = \{s_1^*, \dots, s_{n_y}^*\}$ and \mathcal{R} acts as receiver with input (k, s) . As a result, \mathcal{S} receives nothing and \mathcal{R} receives $b \in \{0, 1\}^{n_y}$, satisfying $b_i = 1$ if and only if s_i^* decrypts to s . Now, we give our framework of mq-RPMT protocol in Figure 7.

Correctness. For all $i \in [n_y]$, if $y_i \in X$, there is an $x_j \in X, j \in [n_x]$ s.t. $y_i = x_j$. In this case, $s_i^* = \text{Decode}_H(D, h(x_j)) = s_j$. Since $s_j = \text{Enc}(k, s)$, we have $\text{Dec}(k, s_j) = s$, which means $b_i = 1$. In the case $y_i \notin X$, if hash functions collide, that is, $h(y_i) = h(x)$ for some $y_i \notin X$, the correctness will be violated. By setting $\sigma = \lambda + \log n_x n_y$, a union bound shows probability of collision is negligible $2^{-\lambda}$. When no collision occurs, from the randomness of OKVS, $s_i^* = \text{Decode}_H(D, h(y_i))$ is a random ciphertext, resulting in s_i^* is not the encryption of s with overwhelming probability. The union bound guarantees that for all $y_i \notin X$, the probability that there exists an s_i^* s.t. $\text{Dec}(k, s_i^*) = s$ is negligible.

We now prove the security properties of our mq-RPMT.

Theorem 3.1 *Assume the encryption scheme \mathcal{E} satisfies single-message multi-ciphertext pseudorandomness. The protocol in Figure 7 securely computes $\mathcal{F}_{\text{mq-rpmt}}$ against semi-honest adversaries in the $\mathcal{F}_{\text{vodm}}$ -hybrid model.*

Proof Due to space limitation, we only sketch here the simulators for the two cases of corrupt \mathcal{S} and corrupt \mathcal{R} , the full proof (via hybrid arguments) is deferred to the full version.

Corrupt sender: To simulate OKVS in Step 3, the simulator computes a random OKVS D by selecting n_x random key-value pairs. Then, the simulator sets $s_i^* := \text{Decode}_H(D, h(y_i))$ and invokes underlying VODM simulator with inputs $(s_1^*, \dots, s_{n_y}^*)$.

Briefly, this simulation is indistinguishable for the following reasons: the single-message multi-ciphertext pseudorandomness of the encryption ensures that value (ciphertext) is indistinguishable from random, and then by the obliviousness of OKVS, D is distributed uniformly.

Corrupt receiver: The simulator for a corrupt \mathcal{R} first obtains b from the ideal mq-RPMT functionality. The only message that needs to be simulated is the VODM functionality in Step 5. The simulator just executes Step 1 honestly and invokes the underlying VODM simulator with inputs (k, s, b) . \square

4 Generic Constructions of Multi-Query RPMT

In this section, we give two generic constructions of mq-RPMT protocol. In the first construction, we use SKE as the encryption scheme and generic 2PC to implement VODM. The advantage is that this scheme only uses OT and symmetric operations. In the second construction, we use PKE and a randomization method to implement the encryption scheme and a leaky version of VODM respectively, which leads to a leaky version of mq-RPMT. However, as observed by KRTW,

Parameters:

- Computational security parameter κ and statistical security parameter λ .
- Input length n .
- A finite group \mathbb{G} .
- Random functions $h_1, h_2, h_3 : \{0, 1\}^* \rightarrow [m']$ and $r : \{0, 1\}^* \rightarrow \{0, 1\}^{d+\lambda}$.
- Parameters $m' = 1.3n$ and $d = 0.5 \log n$, as shown in [19], where d upper bound the size of 2-core of a (m', n) -Cuckoo graph.
- Output length $m = m' + d + \lambda$.

Encode_H({(x₁, y₁), ..., (x_n, y_n)}):

1. Define $l(x) \in \{0, 1\}^{m'}$ to be all zeroes except 1s at positions $h_1(x), h_2(x), h_3(x)$. Here we assume the weight of $l(x)$ is 3. Let $\text{row}(x) := l(x) || r(x)$,

$$M^{(0)} = \begin{bmatrix} l(x_1) \\ \vdots \\ l(x_n) \end{bmatrix} \in \{0, 1\}^{n \times m'}, M^{(1)} = \begin{bmatrix} r(x_1) \\ \vdots \\ r(x_n) \end{bmatrix} \in \{0, 1\}^{n \times (d+\lambda)}$$

and let

$$M = M^{(0)} || M^{(1)} = \begin{bmatrix} \text{row}(x_1) \\ \vdots \\ \text{row}(x_n) \end{bmatrix} \in \{0, 1\}^{n \times m}.$$

2. Initialize empty vectors $D_L \in \mathbb{G}^{m'}$ and $D_R \in \mathbb{G}^{d+\lambda}$, let $D = D_L || D_R$.
3. Initialize stack P .
4. While there is a node $j \in [m']$ such that the set $\{x_i \notin P | j \in \{h_1(x_i), h_2(x_i), h_3(x_i)\}\}$ is a singleton: Let x_i be the element of that singleton, and push x_i onto P .
5. Let $S = \{x_i | x_i \notin P\}$, and let $R \subset [n]$ index the rows of M in S , i.e. $R = \{i | M_{i, h_1(x_i)}^{(0)} = M_{i, h_2(x_i)}^{(0)} = M_{i, h_3(x_i)}^{(0)} = 1 \wedge x_i \in S\}$. Let $\tilde{d} := |R|$ and abort if $\tilde{d} > d$.
6. Let $\tilde{M}^{(1)} \in \{0, 1\}^{\tilde{d} \times (d+\lambda)}$ be the submatrix of $M^{(1)}$ obtained by taking the row indexed by R . Abort if $\tilde{M}^{(1)}$ does not contain an invertible $\tilde{d} \times \tilde{d}$ matrix. Otherwise let \tilde{M}^* be one such matrix and $C \subset [d + \lambda]$ index the corresponding columns of $\tilde{M}^{(1)}$.
7. Let $C' := \{j | i \in R, M_{i, j}^{(0)} = 1\} \cup ([d + \lambda] \setminus C + m')$ and for $i \in C'$ assign $D_i \leftarrow \mathbb{G}$. For $i \in R$, define $y'_i := y_i - (MD^T)_i$ where D_i is assumed to be zero if unassigned.
8. Using Gaussian elimination solve the system $\tilde{M}^* (D_{m'+C_1}, \dots, D_{m'+C_{\tilde{d}}})^T = (y'_{R_1}, \dots, y'_{R_{\tilde{d}}})^T$.
9. While P not empty:
 - (a) pop x_i from P .
 - (b) D_L is undefined in at least one of the positions $h_1(x_i), h_2(x_i), h_3(x_i)$. Set the undefined position(s) so that $\langle \text{row}(x_i), D \rangle = y_i$.
10. Set any empty position in D with a random value from \mathbb{G} .
11. Output D .

Decode_H(D, x):

1. Return $\langle \text{row}(x), D \rangle$.

Figure 3: 3H-GCT++ algorithm

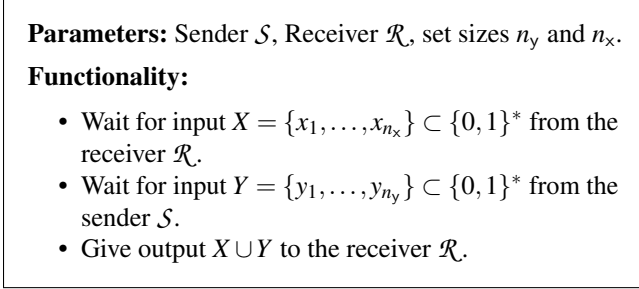


Figure 4: Private Set Union Functionality \mathcal{F}_{psu}

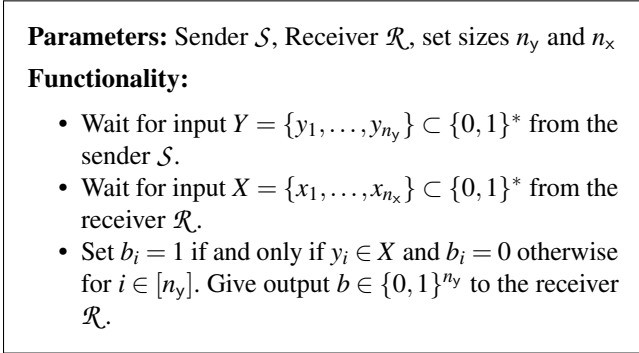


Figure 5: Multi-Query Reverse Private Membership Test Functionality $\mathcal{F}_{\text{mq-rpmt}}$

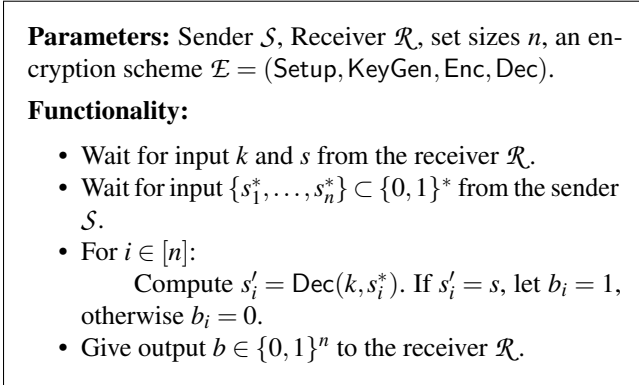


Figure 6: Vector Oblivious Decryption-then-Matching Functionality $\mathcal{F}_{\text{vodm}}$

this leaky version can still be used to construct a secure PSU. Both schemes achieve linear computation and communication complexity.

4.1 Construction from SKE and 2PC

As we noted before, a single-message multi-ciphertext pseudorandom SKE and 2PC are sufficient for constructing mq-RPMT. The correctness and security can be directly derived

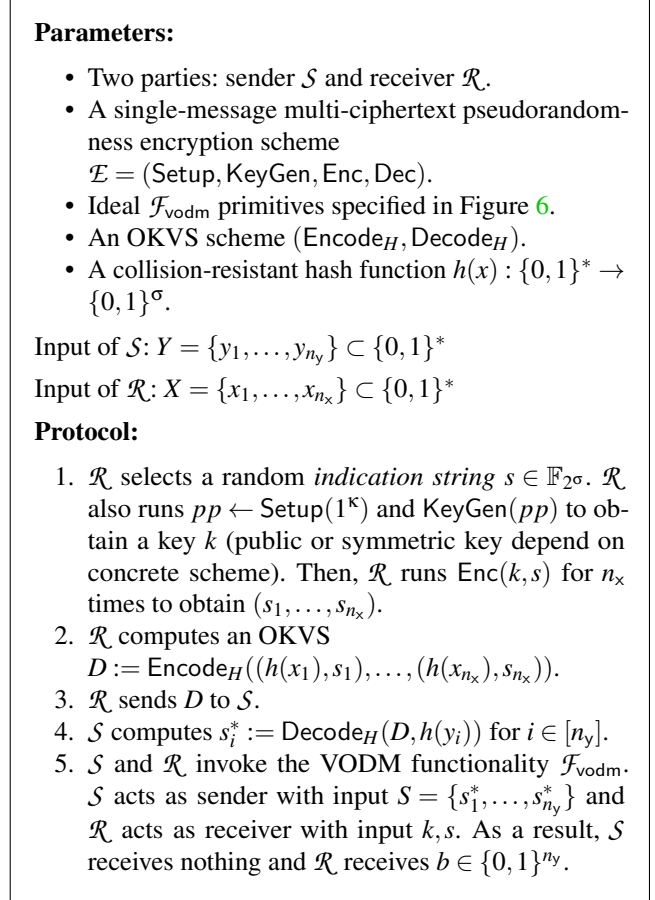


Figure 7: General Construction of mq-RPMT Protocol $\Pi_{\text{mq-rpmt}}$

from the general construction in Section 3.2. It is straightforward to show that PRF-based SKE satisfies the single-message multi-ciphertext pseudorandomness property. We give proof in the full version for completeness.

We use the general 2PC as the implementation of VODM. Formally,

Theorem 4.1 *Taking the PRF-based SKE as the encryption scheme in Figure 7. Assuming that the 2PC implementing VODM is semi-honest secure, then the protocol in Figure 7 securely computes $\mathcal{F}_{\text{mq-rpmt}}$ against semi-honest adversaries.*

This theorem immediately follows from the fact that PRF-based SKE satisfies the single-message multi-ciphertext pseudorandomness property (proved in the full version) and Theorem 3.1.

4.2 Construction from Re-randomizable PKE

Now we consider a specialized way to construct $\mathcal{F}_{\text{vodm}}$. Our main idea is that since the receiver cannot know the randomness used in each ciphertext, as long as the encryption scheme

satisfies the property of rerandomization, the sender can re-randomize all ciphertexts and send the new ciphertexts to the receiver so that the receiver can not obtain additional information by comparing randomness. Note that another problem arises here. The property of re-randomization can only guarantee that for $y \in X$, the receiver is not allowed to learn which one is the sender's element. For $y \notin X$, the ciphertext s_i^* obtained by the sender is related to y , so the plaintext obtained by the receiver is also related to y , which will reveal the information of y . However, as observed by KRTW, in the case of $y \notin X$, we want (in the overall PSU protocol) the receiver to learn y anyway. Fully secure mq-RPMT is actually overkill for PSU, a relaxed version suffices. We define the leaky VODM functionality in Figure 8.

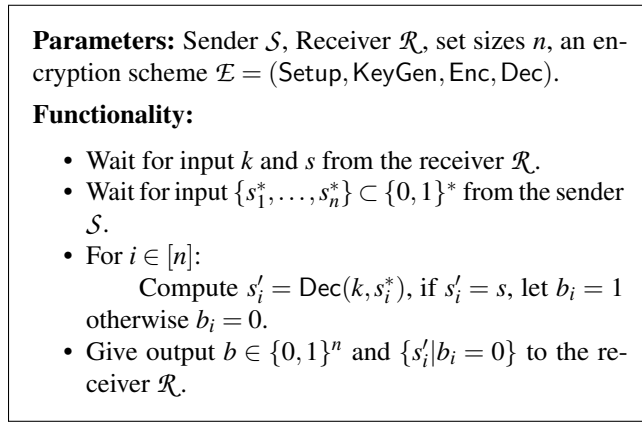


Figure 8: Leaky VODM Functionality $\mathcal{F}_{\text{Vodm}}$

Since the SKE scheme is hard to re-randomize, we consider the use of public-key encryption (PKE) which is easier to re-randomize. We describe our PKE-based leaky VODM protocol in Figure 9.

We now state and prove the security of the above leaky VODM protocol.

Theorem 4.2 *Assume the security of the ReRand-PKE scheme. The protocol in Figure 9 securely computes $\mathcal{F}_{\text{Vodm}}$ against semi-honest adversaries.*

Proof Because the sender does not receive messages in the protocol, we just need to simulate the view of the receiver. We exhibit simulator $\text{Sim}_{\mathcal{R}}$ for simulating corrupt \mathcal{R} .

Corrupt receiver: $\text{Sim}_{\mathcal{R}}(pk, sk, s, b, \{s'_i | b_i = 0\})$ simulates the view of corrupt semi-honest receiver. Note that the only messages that need to be simulated by the simulator are ciphertexts $\{\bar{s}_i\}_{i \in [n]}$.

$\text{Sim}_{\mathcal{R}}$ computes $\bar{s}_i := \text{Enc}(pk, s; r_i)$ if $b_i = 1$ and $\bar{s}_i := \text{Enc}(pk, s'_i; r_i)$ if $b_i = 0$ for $i \in [n]$. $\text{Sim}_{\mathcal{R}}$ appends $\{\bar{s}_i\}_{i \in [n]}$ to the view.

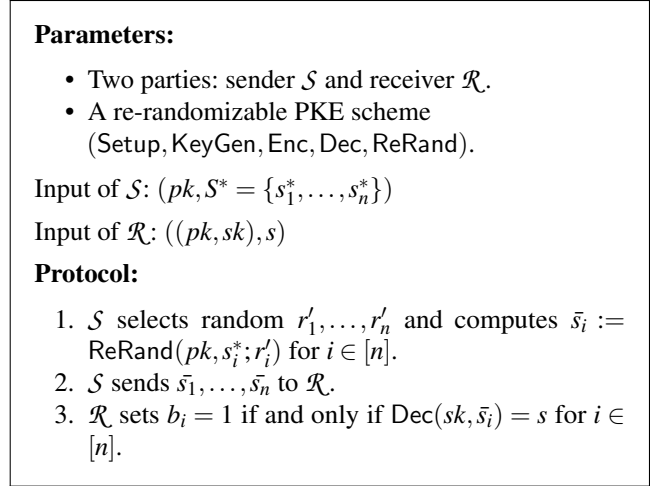


Figure 9: PKE-based Leaky VODM Protocol Π_{Vodm}

The indistinguishability of ReRand-PKE scheme guarantees the view output by $\text{Sim}_{\mathcal{R}}$ is indistinguishable from the real one. \square

Note that the mq-RPMT constructed with the above leaky VODM is also a leaky version. We don't give a specific description of this leaky mq-RPMT. Instead, we use leaky VODM to construct PSU protocol directly and prove its security in the full version.

4.3 Unification with Membership Encryption

We have presented two generic constructions of mq-RPMT protocols from probabilistic SKE and probabilistic PKE respectively. It is intriguing to study if there is a unified way to encompass the two different constructions.

We retrospect the high level idea underlying our mq-RPMT protocol. If privacy is not a concern, reverse membership test can be simply done by having the receiver first create a membership relation R for his set Y , namely $R(y) = 1$ iff $y \in Y$, then having the sender send his elements to the receiver in clear. To make the reverse membership test private, the receiver can "encrypt" his membership relation and send the "encoding" of resulting ciphertexts to the sender. After receiving the "encoding", the sender is able to retrieve the membership encryptions corresponding to his elements. In the sequel, the receiver can fulfill the reverse private membership test by decrypting the ciphertexts in an oblivious manner.

Based on the above discussion, we realize that the right encryption scheme needed in our mq-RPMT protocol is an abstract new notion called *membership encryption (ME)*. Roughly speaking, ME for set X encrypts an element x into a ciphertext, which decrypts to "1" if $x \in X$. We formalize the syntax and security notion of ME in the private-key setting as below.

Definition 4.1 (Membership Encryption) *Membership encryption for set X consists of four polynomial time algorithms satisfying the following properties.*

- $\text{Setup}(1^\kappa)$: on input a security parameter κ , outputs public parameters pp , which include the ciphertext space C .
- $\text{KeyGen}(pp, X)$: on input public parameters pp and $X \subseteq \{0, 1\}^*$, outputs a key k .
- $\text{Enc}(k, x)$: on input a key k and an element $x \in X$, outputs a ciphertext $c \in C$. For uttermost generality, the behavior of Enc on $x \notin X$ is unspecified. Looking ahead, such treatment suffices for the construction of mq-RPMT protocol.
- $\text{Dec}(k, c)$: on input a key k and a ciphertext $c \in C$, outputs “1” indicating c is an encryption of an element x in X and “0” if not.

Correctness. For any $x \in X$, $\forall k \leftarrow \text{KeyGen}(pp, X)$, $\text{Dec}(k, c = \text{Enc}(k, x)) = 1$.

Consistency. For any $x \notin X$, $\Pr[\text{Dec}(k, c) = 0] = 1 - \varepsilon(\kappa)$, where $pp \leftarrow \text{Setup}(1^\kappa)$, $k \leftarrow \text{KeyGen}(pp, X)$, $c \xleftarrow{R} C$. Here, ε is the consistency error, which must be negligible in κ .

Multi-element pseudorandomness. For any n distinct elements $x_1, \dots, x_n \in X$, $\{\text{Enc}(k, x_i)\}_{i \in [n]} \approx_c UC^n$.

The ME notion naturally extends to the public-key setting by letting the KeyGen algorithm generate a keypair (pk, sk) , in which pk is used to encrypt and sk is used to decrypt. We omit the details for its straightforwardness.

We then study the generic construction of ME. Note that the essence of ME is to encrypt element’s membership relation, rather than the element itself. The membership relation can be created by establishing a mapping H from elements to the set under test. Basically, there are two extreme cases of mapping. The first is to select a single indication string s as the characteristic of the set, then map all elements to s , i.e., $H : x_i \rightarrow s$, which we refer to as *lossy mapping*. The second is to select n indication strings s_i as the characteristic of the set, then map elements to distinct indication strings, i.e., $H : x_i \rightarrow s_i$, which we refer to as *injective mapping*. With the above understanding in head, we present various constructions of ME by mixing encryption schemes and membership mapping.

ME from probabilistic SKE and lossy mapping. The construction is as below.

- $\text{Setup}(1^\kappa)$: runs $\text{SKE.Setup}(1^\kappa)$ to generate pp .
- $\text{KeyGen}(pp, X)$: runs $\text{SKE.KeyGen}(pp)$ to sample k_{ske} , picks a random element $s \in M$, where M is the message space of SKE, sets H be a mapping that maps all elements in X to s , outputs $k = (k_{\text{ske}}, H)$
- $\text{Enc}(k, x)$: parses $k = (k_{\text{ske}}, H)$, outputs $c \leftarrow \text{SKE.Enc}(k_{\text{ske}}, H(x))$.
- $\text{Dec}(k, c)$: parses $k = (k_{\text{ske}}, H)$, outputs “1” iff $\text{SKE.Dec}(k_{\text{ske}}, c) = s$.

We note that the above construction can be naturally extended to the public-key case.

Theorem 4.3 *If SKE (resp. PKE) satisfies single-message multi-ciphertext pseudorandomness, then the above ME construction satisfies multi-element pseudorandomness with consistency error $1/|M|$.*

The above ME constructions are exactly the backbones of our generic constructions of mq-RPMT protocol presented in Section 4.1 and 4.2. Since ME requires multi-element pseudorandomness, the use of lossy mapping inherently stipulates that the accompanying encryption schemes are probabilistic. Therefore, in this case the ciphertext expansion is unavoidable. For example, in PRF-based probabilistic SKE, the length of ciphertext is twice that of plaintext. In the design of our mq-RPMT protocol, the value in OKVS is exactly ciphertext. As a consequence, ciphertext expansion incurs overhead to the size of OKVS and thus also the communication cost on the receiver side. For this reason, reducing the ciphertext expansion factor will immediately improve the performance of the overall mq-RPMT protocol.

An important observation is that if we switch to injective mapping, then ME can be built from deterministic encryption schemes satisfying *multi-message multi-ciphertext pseudorandomness*. The constructions are similar as above except the decryption algorithm outputs ‘1’ iff the decryption result falls into the prior-fixed indication string set $S = \{s_i\}_{i \in [n]}$. In instantiation, we take $H : x_i \rightarrow i$ as the membership mapping, which renders efficient membership decryption by testing whether the decryption is less than n .

Formally, we have the following theorem:

Theorem 4.4 *If SKE (resp. PKE) satisfies multi-message multi-ciphertext pseudorandomness, then the ME construction satisfies multi-element pseudorandomness with consistency error $n/|M|$.*

If we instantiate the ME from the PRP-based deterministic SKE and injective mapping, the ciphertext expansion factor is optimal. Therefore, a drop-in replacement to the ME from PRF-based probabilistic SKE and lossy mapping will reduce the size of OKVS in the mq-RPMT protocol by half.

Due to space constraints, we put the description that how to construct mq-RPMT using the language of ME in the full version.

5 Our PSU Protocol

In this section, we describe our PSU construction achieving linear complexity and prove its semi-honest security.

5.1 Generic Construction of PSU Protocols

With mq-RPMT and OT, we can simply combine them to construct a PSU protocol. We give the formal description in

Figure 10.

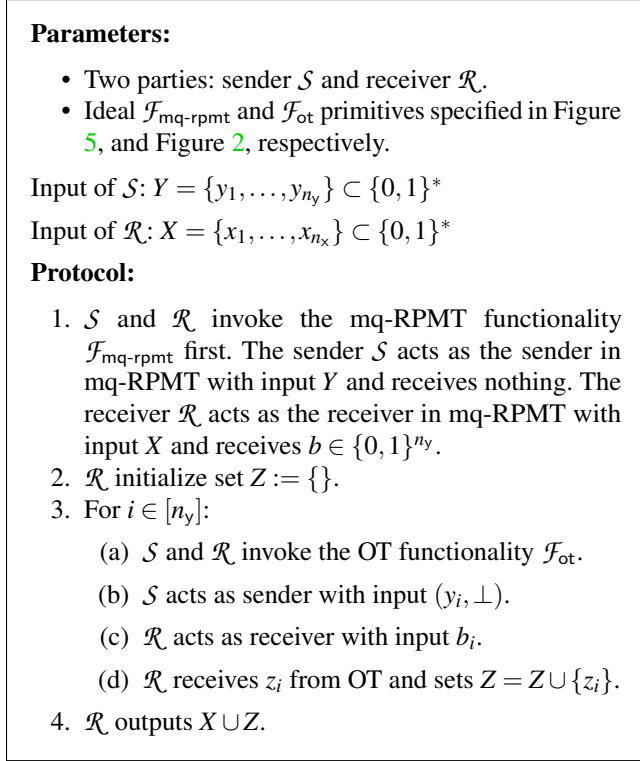


Figure 10: Private Set Union Protocol Π_{psu}

We now state and prove the security properties of the above PSU protocol.

Theorem 5.1 *The protocol in Figure 10 securely computes \mathcal{F}_{psu} against semi-honest adversaries in the $(\mathcal{F}_{\text{mq-rpmt}}, \mathcal{F}_{\text{ot}})$ -hybrid model.*

Proof We exhibit simulators $\text{Sim}_{\mathcal{R}}$ and $\text{Sim}_{\mathcal{S}}$ for simulating corrupt \mathcal{R} and \mathcal{S} respectively, and argue the indistinguishability of the produced transcript from the real execution.

Corrupt Sender: $\text{Sim}_{\mathcal{S}}(Y = \{y_1, \dots, y_{n_y}\})$ simulates the view of corrupt semi-honest sender. It executes as follows:

1. $\text{Sim}_{\mathcal{S}}$ invokes mq-RPMT simulator $\text{Sim}_{\text{mq-rpmt}}^{\mathcal{S}}(Y)$ and appends the output to the view.
2. For $i \in [n_y]$, $\text{Sim}_{\mathcal{S}}$ invokes OT simulator $\text{Sim}_{\text{ot}}^{\mathcal{S}}(y_i, \perp)$ and appends the output to the view.

Now we argue that the view output by $\text{Sim}_{\mathcal{S}}$ is indistinguishable from the real one. This is obtained by the underlying simulators' indistinguishability directly.

Corrupt Receiver: $\text{Sim}_{\mathcal{R}}(X = \{x_1, \dots, x_{n_x}\}, X \cup Y)$ simulates the view of corrupt semi-honest receiver. It executes as follows:

1. $\text{Sim}_{\mathcal{R}}$ defines the set $Z := X \cup Y \setminus X$, i.e. the set of elements that Y “brings to the union”. Next, it uses \perp to pad Z to n_y elements and permutes these elements randomly. Let $Z = \{z_1, \dots, z_{n_y}\}$.
2. $\text{Sim}_{\mathcal{R}}$ sets $b_i = 1$ if and only if $z_i \in X$ for $i \in [n_y]$. Then, it invokes mq-RPMT simulator $\text{Sim}_{\text{mq-rpmt}}^{\mathcal{R}}(X, b)$ and appends the output to the view.
3. For $i \in [n_y]$, $\text{Sim}_{\mathcal{R}}$ invokes OT simulator $\text{Sim}_{\text{ot}}^{\mathcal{R}}(b_i, z_i)$ and appends the output to the view.

Now we argue that the view output by $\text{Sim}_{\mathcal{R}}$ is indistinguishable from the real one. In the simulation, the way \mathcal{R} obtains the elements in $Z = X \setminus Y$ is identical to the real execution. By the underlying simulators' indistinguishability, the simulated view is computationally indistinguishable from the real. \square

5.2 Instantiation of PSU

For our SKE-based construction, we can use a PRP as we mentioned in Section 4.3 to instantiate SKE, which can achieve an optimal ciphertext expansion factor. Since we need to perform the 2PC decryption computation, we use the LowMC [1] as our PRP instantiation to minimize the circuit size. As for generic 2PC, there are two classical methods, e.g. garbled circuit [49] or GMW [22]. The former has a constant number of rounds, while the latter has a lower communication. Since the communication has a greater impact on our scheme, we consider instantiating 2PC by GMW.

For our PKE-based construction, we use the well-known ECC ElGamal [17] scheme as our ReRand-PKE.

5.3 Communication Cost

Now we analyze the communication cost of our two PSU constructions. For the SKE-based construction, we use our ME optimization in Section 4.3.

Let's first analyze the size of decryption circuit in our SKE-based construction: the circuit needs to compute decryption of every $\{s_i^*\}_{i \in [n_y]}$ and compare the result with n_x . If $\text{Dec}(k, s_i^*) < n_x$, it sets $b_i = 1$ and $b_i = 0$ otherwise. The total number of decryption computations is n_y . To compare whether a σ long string is less than n_x , we only need to compute whether the OR of its first $\sigma - \log n_x$ bits are 1, which requires $\sigma - \log n_x - 1$ AND gates (since $a \vee b = \bar{a} \wedge \bar{b}$). The total number of AND gates is $n_y(t + \sigma - \log n_x) = O(tn_y)$, where t is the number of AND gates in a PRP decryption circuit.

Now we are ready to calculate the communication of PSU protocol. Note that the communication of our protocol consists of OKVS, VODM protocol and OT protocol. We analyze their complexity respectively. We use the symbol Φ to represent the communication complexity, and its subscripts represent different components.

Protocol	Communication	$n = n_y = n_x$		
		2^{14}	2^{18}	2^{22}
Frikken [16]	$N + 2n_x N + 4n_y N$	$12288n$	$12288n$	$12288n$
DC [12]	$2\lambda n_x N + 4n_y N$	$172032n$	$172032n$	$172032n$
KRTW [32]	$\beta u(2\rho + \lambda + (u+2)\sigma) + \beta u(\kappa + \sigma)$	$14977n$	$16927n$	$18956n$
GMRSS [18]	$1.27n_y \rho + 3n_x \sigma + (1.27n_y \log n_y + n_y)(\kappa + \sigma)$	$5417n$	$6687n$	$7947n$
JSZDG-R [27]	$\rho(\kappa + 2.18n_x) + 4n_y l_2 + (1.09n_x \log n_x + n_y)(\kappa + \sigma)$	$5757n$	$6931n$	$8105n$
JSZDG-S [27]	$\rho(\kappa + 2.18n_y) + 1.09n_y(ul_2 + \sigma) + 2.18n_y \log n_y(\kappa + \sigma)$	$10640n$	$13140n$	$15658n$
SKE-PSU	$(1.3n_x + d + \lambda)\sigma + \kappa + n_y \sigma + 4n_y(t + \sigma - \log n_x) + n_y(\kappa + \sigma)$	$3768n$	$3810n$	$3853n$
PKE-PSU	$4\kappa(1.3n_x + d + \lambda) + 4\kappa n_y + n_y(\kappa + \sigma)$	$1373n$	$1381n$	$1389n$

Table 2: Theoretical communication costs of PSU protocols (in bits), calculated using computational security $\kappa=128$ and statistical security $\lambda=40$. Ignore costs of base OTs which are independent of input size. N is the size of the public key in Paillier encryption scheme (2048 is used here). β and u are the number of bins and maximum bin size respectively. ρ is the width of OT extension matrix (depends on n and protocol).

- OKVS in both constructions: as we showed in Section 2.5, we use 3H-GCT++ as our OKVS scheme: $\Phi_{\text{okvs}}(n_x) = (1.3n_x + d + \lambda)|c|$, where $|c|$ is the size of ciphertext, $|c| = \lambda + \log n_x n_y$ and 4κ for SKE-based and PKE-based scheme respectively.
- Oblivious decryption:
 - In SKE-based construction: we use $\Phi_{\text{vod}}^{\text{ske}}(n_y, n_x)$ to denote the communication of computing oblivious decryption circuit. As we said in Section 5.2, we use GMW as our 2PC instantiation, the communication consists of *input sharing*, *multiplication gate computation* and *output reconstruction*. In the input sharing phase, the communication is $\kappa + n_y \sigma$ bits, and in the output reconstruction phase, it is n_y bits. Using Beaver triple [4], $4n_y(t + \sigma - \log n_x)$ bits are needed in multiplication phase. So we have $\Phi_{\text{vod}}^{\text{ske}}(n_y, n_x) = \kappa + n_y \sigma + 4n_y(t + \sigma - \log n_x) + n_y$
 - In PKE based construction: the communication of leaky VODM functionality, denoted by $\Phi_{\text{lvodm}}^{\text{pke}}(n_y, n_x) = 4n_y \kappa$
- OT in both constructions: $\Phi_{\text{ot}}(n_y) = n_y(\kappa + \sigma)$.

Let $\Phi_{\text{psu}}^{\text{ske}}(n_y, n_x)$ denote communication of SKE-based construction and let $\Phi_{\text{psu}}^{\text{pke}}(n_y, n_x)$ denote communication of PKE-based construction. The overall communication cost of our PSU protocol is:

$$\Phi_{\text{psu}}^{\text{ske}}(n_y, n_x) = \Phi_{\text{okvs}}(n_x) + \Phi_{\text{vod}}^{\text{ske}}(n_y, n_x) + \Phi_{\text{ot}}(n_y)$$

$$\Phi_{\text{psu}}^{\text{pke}}(n_y, n_x) = \Phi_{\text{okvs}}(n_x) + \Phi_{\text{lvodm}}^{\text{pke}}(n_y, n_x) + \Phi_{\text{ot}}(n_y)$$

5.4 Discussion: Difference between PSI and PSU

Although PSI and PSU are quite similar, as discussed in [32], the techniques they use are different, and building PSU is more challenging than building PSI.

Since the output of PSI is the elements of the receiver’s own set, it is only necessary to test whether each element belongs to the sender’s set (i.e., PMT), and the difficulty of PSU is how to retrieve the elements outside the intersection (i.e., RPMT + OT) without disclosing the intersection. In PSI, PMT can be easily obtained by OPRF: the sender obtains a PRF key k while the receiver obtains $F_k(y)$ on his input y , then the sender computes and sends $\{F_k(x)\}_{x \in X}$ to the receiver. The receiver tests whether $F_k(y) \in \{F_k(X)\}_{x \in X}$ to determine whether $y \in X$. As a result, OPRF is enough for PSI, and all the state-of-the-art PSI protocols [11, 31, 44] follow this paradigm and mainly focus on designing efficient OPRF protocols.

However, the conversion from PMT to RPMT is not trivial, as discussed in [32], this seemingly simple functionality adjustment (PMT \rightarrow RPMT) doesn’t seem to be fixable by a small tweak of PMT. Although OPRF is enough for PSI, this is not the case for PSU. In the state-of-the-art PSU [18, 27], OPRF is only one component, and the design of PSU protocol usually requires the use of a variety of different components, e.g., oblivious switch network functionality, and combine them in a clever method.

5.5 Discussion: the Relationship with Existing PSI/PSU-Related Primitives

Here we also discuss the relationship with existing PSI/PSU-related primitives.

OKVS. Garimella et al. [19] proposed the notion of Oblivious Key-Value Store (OKVS), which is useful in both PSI and PSU. The OKVS is a *data structure* in which a sender has a set of key-value mapping $(\{x_i, y_i\})$ with (pseudo)random y_i ’s, and she wishes to hand that mapping over to a receiver, allowing the receiver to evaluate the mapping on any input but without revealing the keys x_i . Correctness of the data structure must ensure that if the other party evaluates the OKVS on some $q = x_j$ then the result is y_j . Obliviousness here is that the receiver cannot tell what keys x_i ’s are encoded from a given

OKVS. The most compact OKVS that one can think of is a polynomial. The recent excellent works on OKVS [19, 38] make it very efficient to encode a large number of key-value pairs, for example, using 3H-GCT, it takes only about 4.9s to encode 2^{20} key-value pairs.

OTSA. Zhao and Chow [51] proposed a primitives called oblivious transfer for a sparse array (OTSA), which can be used to construct a variant of PSI, i.e. threshold private set intersection (t-PSI). In fact, the OTSA is strictly stronger than OKVS. The OTSA is actually a *protocol* for obviously decoding OKVS, that is, the input of receiver is a set I_r , the input of sender is OKVS $D := \text{Encode}(\{(s_j, e_j)\}_{j \in [n_s]})$, the output of the receiver is $\{\text{Decode}(D, r_j)\}_{j \in [n_r]}$. The main differences between OKVS and OTSA are:

- OTSA enforces the receiver to decode D on limited elements of queries, i.e. I_r , whereas OKVS is simply a data structure that is sent in the clear to the receiver, thus, no limit on the elements of decoding is set.
- In OTSA, the receiver does not know the correspondence between r_j and $\text{Decode}(D, r_j)$ (i.e. sender indices privacy), while in OKVS, the receiver directly knows the relationship between r_j and $\text{Decode}(D, r_j)$.

These limitations have a significant impact on their performance, for example, the experiment in [51] showed that their most efficient OTSA protocol takes about 400s for input size $n = 2^{10}$. It is enough for our construction to use simpler and more efficient OKVS instead of OTSA.

OVDM. In our PSU construction, we proposed a new primitive called oblivious vector decryption-then-matching (OVDM), which is also a *protocol* aiming to decrypt a vector of ciphertexts obviously and then match the decrypted ciphertext to a given string. The significant differences between OVDM and OTSA are:

- OTSA is the protocol for decoding an *OKVS*, while OVDM is the protocol for decrypting an *encryption scheme*.
- OTSA allows the party providing the decoding material (i.e. I_r) to obtain the decoding result (since the decoding algorithm is written as $\text{Decode}(D, r_j)$, D can be regarded as a "key" in some sense), while OVDM allows the party providing the key to obtain the decryption result.
- The output of OVDM is only 1 bit information of plaintext, i.e. whether the plaintext is equal to a string input by the receiver.
- The order of the decryption results output by OVDM is the same as the order of the ciphertext input by the sender, while OTSA does not preserve this order (i.e. sender indices privacy).

Due to the above differences, the ideas for constructing OTSA and OVDM are different. Our OVDM is more efficient than OTSA because we only need PKE to meet the Re-rand property, while OTSA requires more complex homomorphic

PKE.

One may wonder whether the construction of OVDM depends on the particular OKVS construction. We clarify that OVDM and OKVS are two different notions of different usages. We use the combination of OKVS and OVDM to construct mqRPMT, as shown in Section 3. Any OKVS instantiation that meets Randomness can be used for our mqRPMT construction. The only connection between OKVS and OVDM is that they share the same encryption scheme, that is, the value encoded by OKVS is the ciphertext of the encryption scheme, and the sender takes the ciphertext decoded from OKVS as her OVDM input. Since decryption is required, the construction of OVDM is related to the selection of encryption schemes (therefore, we classify our schemes as SKE-based and PKE-based).

6 Implementation and Performance

Recall that we have presented two variants of our protocol. In this section, we will refer to them as:

- SKE-PSU: PSU protocol with SKE-based mq-RPMT, where SKE and VODM are instantiated with PRP and GMW [22] respectively.
- PKE-PSU: PSU protocol with PKE-based mq-RPMT, where ReRand-PKE is instantiated with ECC ElGamal encryption scheme.

The OKVS instantiation of both schemes uses the 3H-GCT++ in Figure 3. We focus on the case where $n_y = n_x = n$, i.e., both parties have equal-size sets.

6.1 Theoretical Analysis of Communication

In Table 2, we show the theoretical communication complexity of our protocol compared with the Frikken protocol [16], the DC protocol [12], the KRTW protocol [32], the GMRSS protocol [18] and the JSZDG protocol [27] (note that [27] proposed two protocols, i.e. JSZDG-R and JSZDG-S, which focus on balanced and unbalanced setting, respectively) in the semi-honest setting. Empirical measurements of such real-world costs are given later in Table 3.

For set sizes in the range 2^{14} to 2^{22} , our PKE-PSU variant has the least communication of any of the protocols we consider: up to an $8.8\times$ improvement of Frikken, $125\times$ improvement of DC, $10.9 - 13.6\times$ improvement of KRTW, $3.9 - 5.7\times$ improvement of GMRSS, and $4.2 - 11.3\times$ improvement of JSZDG. It means that our scheme has great advantages in low bandwidth scenarios.

For our SKE-based protocol, as mentioned in Section 5.2, we use LowMC [1] to minimize the number of AND gates. Though the communication of our SKE-PSU protocol is about $3\times$ higher than PKE-PSU, it is still lower than all previous schemes.

6.2 Experimental Setup

We run our experiments on a single Intel Core i9-9900K with 3.6GHz and 128GB RAM. We simulate the network connection using Linux `tc` command. To better meet the potential deployment requirements, we use Netty² to maintain the communication channel. And we use Protocol Buffers³ for data (de-)serialization. Refer to the full version for details of Netty and Protocol Buffers.

6.3 Implementation Details

Existing PSU implementations are under different MPC frameworks and different experimental settings. For example, the [32] implementation is under 128-bit element length while the [18] implementation is under 64-bit element length. Also, the [27, 32] implementation supports multi-thread execution, while the [18] implementation does not. Further, the [18] and [27] implementation heavily relies on 1-out-of-2 Oblivious Transfer (OT). Introducing recent silent OT technique may further reduce its communication cost [6, 48]. However, existing efficient silent OT implementation [48] is available in `emp-toolkit` [47]. Combining these implementations rely on relatively heavy source code modification works.

After carefully studying existing open-source codes, we fully re-implement state-of-the-art PSU protocols [18, 27, 32] and their underlying basic protocols using Java, including base OT [36], OT extension [2], silent OT [48], the specific OPRF variant [31], and GCT data structures.

We choose Java as our primary programming language for the following reasons. First, recent advances in MPC make this attractive data security technique from theory into practical usage. Introducing big data frameworks into MPC would further increase its efficiency and integrate MPC with existing data pipelines [3]. Current widely adopted big data analytical engines, for example, Hadoop and Spark, are built upon Java or JVM-based programming languages. We hope our implementation can help developers from the big data community leverage and deploy MPC in a more scaling manner. Second, one may think that Java is much slower than C/C++. It is shown that although there is some performance gap, most basic operations in Java and C/C++ have similar performances⁴.

For operations that have a huge efficiency gap between Java and C/C++, we use the Java Native Interface (JNI) technique to invoke C/C++ libraries. The details can be found in the full version.

We note that our implementations support multi-thread executions for all the PSU schemes, including [18], achieved by using Java `'Stream.parallel()'`. In our

²<https://netty.io/>

³<https://developers.google.com/protocol-buffers>

⁴Our tests show that on Macbook Pro 2019, Java needs 0.095us for one AES operation, while C/C++ under AES instruction needs 0.071us. This is because Java would automatically use AES instruction if it detects that the current operating system supports it.

experiments, we limit the number of threads during the protocol execution by setting the JVM parameter `'java.util.concurrent.ForkJoinPool.common.parallelism'`, and submit all parallel executions into that common thread pool. In the single-thread setting, we let all procedures run in the main thread instead of simply setting the number of threads to be one under the multi-thread setting, thus avoiding additional costs for creating and destroying sub-threads. Our performance reports show that we obtained improved performance results for the [18] PSU scheme.

Although most operations in Java and C/C++ have similar performances, there are some operations in which Java operates much slower than C/C++. For example, our JSZDG performance results (See Table 3) are about 3 times slower than the report shown in the original work [27]. We carefully analyze our implementation and find that the gap is from its underlying batch OPRF proposed by Chase and Miao [11]. Briefly speaking, this batched OPRF needs to map each element into a long pseudo-random byte array via a PRF and then convert that to be an integer array as coordinates in the random encoding matrix. In C/C++, the transformation can be done by simply changing the pointer type from `uint8_t*` to `uint32_t*` with almost no additional cost. However, such an operation is not supported in Java due to the memory protection mechanism. One has to explicitly convert `byte[]` to `int[]`, involving dramatic costs. In addition, the type conversion operation cost is, unfortunately, lower than JNI invoking. Introducing JNI in this operation leads to even more costs. How to efficiently implement the batch OPRF proposed by Chase and Miao [11] in memory-safe programming language as in C/C++ remains an open problem in the implementation. We emphasize that designing a unified framework for all PSU protocols while compatible with widely adopted big data analytical engines under C/C++ would further lead to better performance results. We hope that our implementation can be a starting point. Our complete implementation is available on GitHub: <http://github.com/alibaba-edu/mpc4j>.

6.4 Experimental Details

The SKE-PSU protocol is instantiated with the LowMC encryption scheme [1] where the block size and the key length are both 128 bits, and the number of Sboxes is $m = 10$ (i.e., the SboxLayer is a 10-folded parallel application of the basic 3-bit Sbox on the first 30 bits of the state, and for the remaining 88 bits, the SboxLayer is the identity). The concrete parameters in LowMC are from the Mobile PSI implementations provided by Kales et al. [28]⁵. We use the improved inverse of the SBoxLayer provided by Liu et al. [34] and follow the SBoxLayer implementation idea by Kales et al. [28] to implement the (non-2PC) decryption procedure. The

⁵https://github.com/contact-discovery/mobile_psi_cpp/blob/master/droidCrypto/lowmc/lowmc_128_128_20.c

n	Protocol	Comm. (MB)					Running time (s)															
		\mathcal{R}		S		total	LAN				1Gbps				100Mbps				10Mbps			
		setup	online	setup	online		$T=1$	$T=8$	$T=1$	$T=8$	$T=1$	$T=8$	$T=1$	$T=8$	$T=1$	$T=8$	$T=1$	$T=8$				
						setup													online	setup	online	setup
2^{14}	KRTW	0.02	4.17	0.01	29.63	33.8	0.07	3.5	0.03	1.07	0.49	16.13	0.37	14.06	0.83	27.36	0.72	24.66	0.81	55.9	0.73	55.32
	GMRSS	0.02	5.89	0.02	7.96	13.85	0.1	1.01	0.04	0.42	0.66	1.96	0.46	1.28	1	3.53	0.91	2.97	1.06	14.44	0.93	13.97
	JSZDG-R	0.01	4.65	0.01	5.63	10.28	0.07	1.81	0.02	0.52	0.27	2.65	0.23	1.34	0.49	4.19	0.41	2.66	0.45	12.08	0.37	10.63
	SKE-PSU	0.01	3.16	0	3.36	6.52	0.03	0.65	0.02	0.29	0.12	6.76	0.11	6.48	0.21	12.66	0.19	12.09	0.2	15.62	0.19	15.59
	PKE-PSU	0.01	1.16	0	1.59	2.75	4.6	2.37	4.58	1.07	4.78	2.63	4.75	1.34	4.92	3.02	4.9	1.77	4.99	4.43	4.91	3.79
	PKE-PSU*	0.01	2.16	0	2.9	5.05	4.6	1.96	4.6	0.59	4.75	2.36	4.76	1	4.95	2.76	4.91	1.54	4.92	5.72	4.93	5.31
2^{16}	KRTW	0.02	17.64	0.01	122.05	139.69	0.07	12.57	0.03	3.76	0.46	26.27	0.39	20.96	0.82	40.09	0.73	36.3	0.81	163.48	0.75	161.63
	GMRSS	0.02	25.95	0.02	34.11	60.06	0.11	4.79	0.04	1.95	0.64	6.61	0.48	4.25	1.11	12.67	0.92	9.78	1.04	60.75	0.94	57.5
	JSZDG-R	0.01	20.75	0.01	24.74	45.49	0.07	7.5	0.02	2.25	0.3	9.29	0.2	4.45	0.44	13.78	0.4	8.58	0.47	49.41	0.42	44.58
	SKE-PSU	0.01	12.61	0	13.41	26.03	0.04	2.66	0.02	1.15	0.13	8.66	0.11	7.32	0.2	15.84	0.19	14.39	0.2	31.79	0.19	30.98
	PKE-PSU	0.01	4.62	0	6.37	10.99	4.62	9.75	4.59	4.39	4.13	10.21	4.76	5.22	4.9	10.94	4.91	5.83	5.01	16.38	4.92	13.61
	PKE-PSU*	0.01	8.63	0	11.57	20.19	4.57	7.96	4.6	2.58	4.76	8.68	4.77	3.37	4.93	9.94	4.91	4.65	4.94	21.46	4.93	19.67
2^{18}	KRTW	0.02	69.29	0.01	562.76	632.05	0.08	63.02	0.03	17.67	0.52	85.56	0.39	45.31	0.76	111.14	0.71	113.83	0.84	660.33	0.74	664.93
	GMRSS	0.02	113.7	0.02	145.11	258.81	0.13	20.74	0.03	9.8	0.58	28.62	0.55	16.63	1.09	49.68	0.93	38.82	1.03	251.84	0.97	243.63
	JSZDG-R	0.01	92.67	0.01	107.89	200.56	0.07	41.15	0.03	10.71	0.25	43.17	0.21	16.84	0.42	64.06	0.4	33.8	0.53	221.27	0.39	191.2
	SKE-PSU	0.01	50.34	0	53.51	103.85	0.04	10.78	0.02	4.88	0.12	17.83	0.1	12.32	0.2	28.38	0.18	22.54	0.21	98.96	0.19	95.72
	PKE-PSU	0.01	18.5	0	25.45	43.95	4.6	41.5	4.59	19.82	4.79	42.37	4.75	20.97	4.92	44.8	4.91	23.38	4.92	66.68	4.9	54.39
	PKE-PSU*	0.01	34.5	0	46.26	80.76	4.61	34.63	4.58	12.26	4.78	37.1	4.75	13.99	4.92	40.62	4.92	18.45	4.91	85.31	4.92	79.22
2^{20}	KRTW	0.02	300.14	0.01	2305.8	2605.95	0.11	245.37	0.04	67.97	0.52	281.96	0.38	120.35	0.82	363.95	0.74	361.12	0.84	2643.84	0.75	2638.05
	GMRSS	0.02	493.2	0.02	615.9	1109.1	0.11	100.48	0.04	48.53	0.62	119.98	0.51	75.76	1.11	207.83	0.95	164.25	1.09	1074.33	0.95	1030.3
	JSZDG-R	0.01	405.53	0.01	467.26	872.79	0.08	173.07	0.04	54.41	0.48	184.63	0.2	73.28	0.47	266.51	0.73	146.13	0.47	941.5	0.72	825.16
	SKE-PSU	0.01	200.88	0	213.55	414.43	0.05	44.73	0.03	22.78	0.13	59.65	0.11	35.71	0.2	86.11	0.2	65.18	0.21	378.57	0.4	369.24
	PKE-PSU	0.01	74	0	101.8	175.8	4.65	168.79	4.6	79.95	4.78	169.18	4.79	86.49	4.97	179.58	4.94	96.32	4.97	269.32	4.87	216.19
	PKE-PSU*	0.01	138	0	185	323	4.64	144.24	4.58	50.56	4.75	146.41	4.74	60.5	4.9	161.26	5	76.33	4.99	345	4.9	313.37

Table 3: Communication cost (in MB) and running time (in seconds) comparing our protocols to KRTW GMRSS, and JSZDG-R. The LAN network has 10 Gbps bandwidth and 0.2 ms RTT latency. Communication cost of S/\mathcal{R} indicates the outgoing communication from S/\mathcal{R} to the other party. The best protocol within a setting is marked in blue.

underlying OKVs for our PSU protocols are instantiated with our 3H-GCT++ in Figure 3.

Since both [18] and [27] protocols rely heavily on OSN [35] and involve a large number of OT. We further introduce Silent OT [6, 48] in the GMRSS and JSZDG schemes. See details in the full version.

In SKE-PSU, we assume a commonly used setting where Boolean multiplication triples are pre-computed offline and stored locally in a temporary file. This follows real scenarios where Boolean multiplication triples are pre-generated by parties themselves or with the help of a Trusted-Third Party under the Trusted Dealer model. For completeness, we give the costs of triple generation in the full version.

In PKE-PSU, the ReRand-PKE is instantiated with the ECC ElGamal encryption scheme under the curve SecP256K1. We found an interesting point in the implementation of PKE-PSU: In elliptic-curve-based cryptography, point compression is a standard trick, which can roughly reduce the representation of an EC point by half. The cost of this trick is that one has to perform point decompression in the future, which is typically considered to be cheap. Somewhat surprisingly, it turns out that point decompression is very costly. According to existing implementations provided in MCL and OpenSSL libraries, point decompression is as expensive as point exponentiation. Due to this fact, we prefer to use standard point representation for better efficiency when bandwidth is not of first priority. In the implementation, we use **PKE-PSU*** to represent the version that does not perform point compression.

The simulated network settings include typical LAN

(10Gbps bandwidth and 0.02ms RTT latency) and WAN (including 1Gbps with 40ms latency, 100Mbps bandwidth with 80ms latency). In our KRTW implementation, we follow the pipelining optimization shown in [32] with 2^8 pipelining size when the receiver sends polynomials to the sender. In our PKE-PSU, we also leverage the pipelining optimization with the same 2^8 pipelining size when the sender sends ReRand outputs to the receiver.

We divide all protocols into two phases: the one-time setup phase and the online phase. As the name suggests, the one-time setup phase does necessary operations before actual protocol execution, including key distribution, base OT execution, and the one-time setup phase for Ferret OT [48]. The online phase does subsequent protocol executions. Note that in our PKE-PSU, the receiver can send the public key to the sender in the one-time setup phase, and all fixed-point pre-computations related to the public key can also be done in that phase. We emphasize that fixed-point precomputations only need to be performed once, regardless of the number of subsequent protocol executions.

Since the JSZDG-S scheme [27] focus on unbalanced setting and its performance is about $2\times$ worse than the JSZDG-R scheme, we only compare our schemes with JSZDG-R here. Detailed comparisons for set sizes 2^{14} , 2^{16} , 2^{18} , 2^{20} and controlled network configurations are shown in Table 3. To be more intuitive, we show the variation of the running time with the bandwidth in different setting in Figure 11.

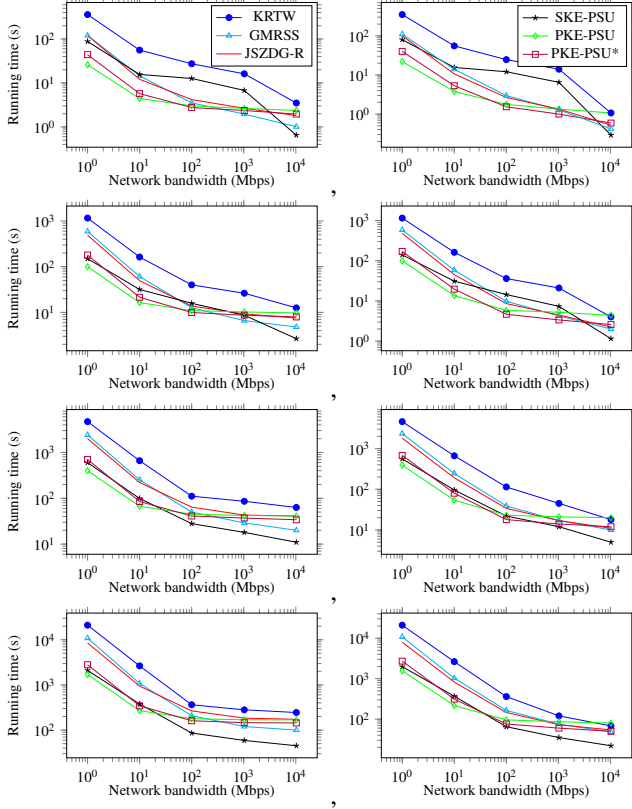


Figure 11: Decline of running time (in seconds) on increasing network bandwidth for our protocols compared with KRTW, GMRSS and JSZDG-R. Both x and y -axis are in log scale. The four figures on the left correspond to $T = 1$ and the right correspond to $T = 8$. The corresponding set sizes from the first row to the last row are $n = 2^{14}, 2^{16}, 2^{18}, 2^{20}$ respectively.

6.5 Performance Evaluation

Communication improvement. As shown in Table 3, our PKE-PSU protocol has the lowest communication among all protocols, which is $12.3 - 14.8\times$ lower than KRTW, $5.1 - 6.3\times$ lower than GMRSS and $3.7 - 5\times$ lower than JSZDG-R. The communication of PKE-PSU* is about $2\times$ higher than that of PKE-PSU, which is due to the absence of point compression. The communication of our SKE-PSU is about $2.5\times$ higher than that of PKE-PSU. Nevertheless, all our schemes have lower communication than that of KRTW, GMRSS and JSZDG-R schemes. Since the communication costs of all our protocols are linear with the parties' set sizes, while the communication costs of the other protocols are not. The larger the parties' set sizes are, the larger the communication cost ratios are.

Computation improvement. As shown in Table 3 and Figure 11, our SKE-PSU performs best when the set size and the bandwidth are large. For example, for $n = 2^{20}$ with $T = 1$ thread in LAN setting, SKE-PSU requires 44.73 seconds,

achieving a $5.5\times$ improvement over KRTW, a $2.2\times$ improvement over GMRSS, and a factor of $3.9\times$ improvement over JSZDG-R.

Our PKE-PSU and PKE-PSU* could be seen as a trade-off between communication and computation. Both schemes perform better in lower bandwidth. Our PKE-PSU scheme is the fastest one under 10Mbps, which is due to its lowest communication, e.g., for $n = 2^{20}$, PKE-PSU requires 216.19 seconds with $T = 8$ threads, while KRTW requires 2638.05 seconds, a $12.2\times$ improvement, GMRSS requires 1030.3 seconds, a $4.8\times$ improvement, and JSZDG-R requires 825.16 seconds, a $3.8\times$ improvement. Our PKE-PSU* performs better in medium bandwidth (100Mbps and 1Gbps). For example, for $n = 2^{18}$ with $T = 8$ threads in 100Mbps, PKE-PSU* requires 18.45 seconds, while KRTW requires 113.83 seconds, a $6.2\times$ improvement, GMRSS requires 38.82 seconds, a $2.1\times$ improvement, and JSZDG-R requires 33.8 seconds, a $1.8\times$ improvement. We also noticed that the performance of PKE-PSU* improved significantly (about $3\times$ speedup) in the case of multithreading because of its heavy computation cost.

6.6 Applications

We further gave the experiment results of two PSU applications introduced in Section 1, namely IP blacklist aggregation and Private ID. Due to space limitations, the detailed experiment is shown in the full version.

Acknowledgement

We are grateful for the helpful comments from the anonymous reviewers. Cong Zhang and Dongdai Lin are supported by the National Key Research and Development Program of China (No. 2020YFB1805402) and the National Natural Science Foundation of China (Grants No. 61872359 and No. 61936008). Yu Chen and Min Zhang are supported by the National Key Research and Development Program of China (Grant No. 2021YFA1000600) and the National Natural Science Foundation of China (Grant No. 62272269).

References

- [1] Martin R. Albrecht, Christian Rechberger, Thomas Schneider, Tyge Tiessen, and Michael Zohner. Ciphers for MPC and FHE. In *EUROCRYPT 2015*, 2015.
- [2] Gilad Asharov, Yehuda Lindell, Thomas Schneider, and Michael Zohner. More efficient oblivious transfer and extensions for faster secure computation. In *CCS 2013*, 2013.
- [3] Saikrishna Badrinarayanan, Ranjit Kumaresan, Mihai Christodorescu, Vinjith Nagaraja, Karan Patel, Srinivasan Raghuraman, Peter Rindal, Wei Sun, and Minghua Xu. A plug-n-play framework for scaling private set

- intersection to billion-sized sets. Cryptology ePrint Archive, Paper 2022/294, 2022. <https://eprint.iacr.org/2022/294>.
- [4] Donald Beaver. Efficient multiparty protocols using circuit randomization. In *CRYPTO 1991*, 1991.
- [5] Marina Blanton and Everaldo Aguiar. Private and oblivious set and multiset operations. In *ASIACCS 2012*, 2012.
- [6] Elette Boyle, Geoffroy Couteau, Niv Gilboa, Yuval Ishai, Lisa Kohl, Peter Rindal, and Peter Scholl. Efficient two-round OT extension and silent non-interactive secure computation. In *CCS 2019*, 2019.
- [7] Justin Brickell and Vitaly Shmatikov. Privacy-preserving graph algorithms in the semi-honest model. In *ASIACRYPT 2005*, 2005.
- [8] Prasad Buddharapu, Andrew Knox, Payman Mohassel, Shubho Sengupta, Erik Taubeneck, and Vlad Vlaskin. Private matching for compute. Cryptology ePrint Archive, Paper 2020/599, 2020. <https://eprint.iacr.org/2020/599>.
- [9] M. Burkhart and Xenofontas Dimitropoulos. Fast private set operations with sepia. 2012.
- [10] Ran Canetti, Omer Paneth, Dimitrios Papadopoulos, and Nikos Triandopoulos. Verifiable set operations over outsourced databases. In *PKC*, 2014.
- [11] Melissa Chase and Peihan Miao. Private set intersection in the internet setting from lightweight oblivious PRF. In *CRYPTO 2020*, 2020.
- [12] Alex Davidson and Carlos Cid. An efficient toolkit for computing private set operations. In *ACISP 2017*, 2017.
- [13] Changyu Dong, Liqun Chen, and Zikai Wen. When private set intersection meets big data: an efficient and scalable protocol. In *CCS 2013*, 2013.
- [14] Brett Hemenway Falk, Daniel Noble, and Rafail Ostrovsky. Private set intersection with linear communication from general assumptions. In *WPES@CCS 2019*, 2019.
- [15] Michael J. Freedman, Kobbi Nissim, and Benny Pinkas. Efficient private matching and set intersection. In *EUROCRYPT 2004*, 2004.
- [16] Keith B. Frikken. Privacy-preserving set union. In *ACNS 2007*, 2007.
- [17] Taher El Gamal. A public key cryptosystem and a signature scheme based on discrete logarithms. *IEEE Trans. Inf. Theory*, 31(4):469–472, 1985.
- [18] Gayathri Garimella, Payman Mohassel, Mike Rosulek, Saeed Sadeghian, and Jaspal Singh. Private set operations from oblivious switching. In *PKC 2021*, 2021.
- [19] Gayathri Garimella, Benny Pinkas, Mike Rosulek, Ni Trieu, and Avishay Yanai. Oblivious key-value stores and amplification for private set intersection. In *CRYPTO 2021*, 2021.
- [20] Satrajit Ghosh and Tobias Nilges. An algebraic approach to maliciously secure private set intersection. In *EUROCRYPT 2019*, 2019.
- [21] Oded Goldreich. *The Foundations of Cryptography - Volume 2: Basic Applications*. Cambridge University Press, 2004.
- [22] Oded Goldreich, Silvio Micali, and Avi Wigderson. How to play any mental game or A completeness theorem for protocols with honest majority. In *STOC 1987*, 1987.
- [23] Carmit Hazay and Kobbi Nissim. Efficient set operations in the presence of malicious adversaries. In *PKC 2010*, 2010.
- [24] Kyle Hogan, Noah Luther, Nabil Shear, Emily Shen, David Stott, Sophia Yakoubov, and Arkady Yerukhimovich. Secure multiparty computation for cooperative cyber risk assessment. In *SecDev 2016*, 2016.
- [25] Jeongdae Hong, Jung Woo Kim, Jihye Kim, Kunsoo Park, and Jung Hee Cheon. Constant-round privacy preserving multiset union. Cryptology ePrint Archive, Report 2011/138, 2011. <https://ia.cr/2011/138>.
- [26] Yuval Ishai, Joe Kilian, Kobbi Nissim, and Erez Petrank. Extending oblivious transfers efficiently. In *CRYPTO 2003*, 2003.
- [27] Yanxue Jia, Shi-Feng Sun, Hong-Sheng Zhou, Jiajun Du, and Dawu Gu. Shuffle-based private set union: Faster and more secure. In *USENIX Security 22*, 2022.
- [28] Daniel Kales, Christian Rechberger, Thomas Schneider, Matthias Senker, and Christian Weinert. Mobile private contact discovery at scale. In *USENIX Security 2019*, 2019.
- [29] Lea Kissner and Dawn Xiaodong Song. Privacy-preserving set operations. In *CRYPTO 2005*, 2005.
- [30] Vladimir Kolesnikov and Ranjit Kumaresan. Improved OT extension for transferring short secrets. In *CRYPTO 2013*, 2013.
- [31] Vladimir Kolesnikov, Ranjit Kumaresan, Mike Rosulek, and Ni Trieu. Efficient batched oblivious PRF with applications to private set intersection. In *CCS 2016*, 2016.
- [32] Vladimir Kolesnikov, Mike Rosulek, Ni Trieu, and Xiao Wang. Scalable private set union from symmetric-key techniques. In *ASIACRYPT*, 2019.
- [33] Arjen K. Lenstra and Tim Voss. Information security risk assessment, aggregation, and mitigation. In *ACISP 2004*, 2004.
- [34] Fukang Liu, Takanori Isobe, and Willi Meier. Cryptanalysis of full lowmc and lowmc-m with algebraic techniques. In *CRYPTO 2021*, 2021.
- [35] Payman Mohassel and Seyed Saeed Sadeghian. How to hide circuits in MPC an efficient framework for private function evaluation. In *EUROCRYPT 2013*, 2013.

- [36] Moni Naor and Benny Pinkas. Efficient oblivious transfer protocols. In *Proceedings of the Twelfth Annual Symposium on Discrete Algorithms*, 2001.
- [37] Benny Pinkas, Mike Rosulek, Ni Trieu, and Avishay Yanai. Spot-light: Lightweight private set intersection from sparse OT extension. In *CRYPTO 2019*, 2019.
- [38] Benny Pinkas, Mike Rosulek, Ni Trieu, and Avishay Yanai. PSI from paxos: Fast, malicious private set intersection. In *EUROCRYPT 2020*, 2020.
- [39] Benny Pinkas, Thomas Schneider, Gil Segev, and Michael Zohner. Phasing: Private set intersection using permutation-based hashing. In *USENIX Security 2015*, 2015.
- [40] Benny Pinkas, Thomas Schneider, and Michael Zohner. Faster private set intersection based on OT extension. In *USENIX Security*, 2014.
- [41] Benny Pinkas, Thomas Schneider, and Michael Zohner. Scalable private set intersection based on OT extension. *ACM Trans. Priv. Secur.*, 21(2):7:1–7:35, 2018.
- [42] Michael O. Rabin. How to exchange secrets with oblivious transfer. *IACR Cryptol. ePrint Arch.*, 2005:187, 2005.
- [43] Sivaramakrishnan Ramanathan, Jelena Mirkovic, and Minlan Yu. BLAG: improving the accuracy of blacklists. In *NDSS*, 2020.
- [44] Peter Rindal and Phillipp Schoppmann. VOLE-PSI: fast OPRF and circuit-psi from vector-ole. In *EUROCRYPT 2021*, 2021.
- [45] Jae Hong Seo, Jung Hee Cheon, and Jonathan Katz. Constant-round multi-party private set union using reversed laurent series. In *PKC 2012*, 2012.
- [46] Katsunari Shishido and Atsuko Miyaji. Efficient and quasi-accurate multiparty private set union. In *SMART-COMP 2018*, 2018.
- [47] Xiao Wang, Alex J. Malozemoff, and Jonathan Katz. EMP-toolkit: Efficient MultiParty computation toolkit. <https://github.com/emp-toolkit>, 2016.
- [48] Kang Yang, Chenkai Weng, Xiao Lan, Jiang Zhang, and Xiao Wang. Ferret: Fast extension for correlated OT with small communication. In *CCS 2020*, 2020.
- [49] Andrew Chi-Chih Yao. How to generate and exchange secrets (extended abstract). In *FOCS*, 1986.
- [50] Cong Zhang, Yu Chen, Weiran Liu, Min Zhang, and Dongdai Lin. Optimal private set union from multi-query reverse private membership test. *Cryptology ePrint Archive*, Paper 2022/358, 2022. <https://eprint.iacr.org/2022/358>.
- [51] Yongjun Zhao and Sherman S. M. Chow. Are you the one to share? secret transfer with access structure. *Cryptology ePrint Archive*, Paper 2015/929, 2015. <https://eprint.iacr.org/2015/929>.