

Fast Unbalanced Private Set Union from Fully Homomorphic Encryption

Binbin Tu¹, Yu Chen¹, Qi Liu¹, and Cong Zhang^{2,3}

¹ School of Cyber Science and Technology, Shandong University

² State Key Laboratory of Information Security, Institute of Information Engineering, Chinese Academy of Sciences

³ School of Cyber Security, University of Chinese Academy of Sciences

{tubinbin, liuqicst}@mail.sdu.edu.cn; yuchen@sdu.edu.cn; zhangcong@iie.ac.cn

Abstract. Private set union (PSU) allows two parties to compute the union of their sets without revealing anything else. It has been widely used in various applications. While several computationally efficient PSU protocols have been developed for the balanced case, they have a potential limitation in their communication complexity, which grows (super)-linearly with the size of the larger set. This poses a challenge when performing PSU in the unbalanced setting, where one party is a constrained device holding a small set, and another is a service provider holding a large set.

In this work, we propose a generic construction of unbalanced PSU from leveled fully homomorphic encryption and a newly introduced protocol called permuted matrix private equality test. By instantiating the generic construction, we obtain two unbalanced PSU protocols whose communication complexity is linear in the size of the smaller set, and logarithmic in the larger set.

We implement our protocols. Experiments demonstrate that our protocols outperform all previous protocols in the unbalanced setting. The larger difference between the sizes of two sets, the better our protocols perform. For input sets of sizes 2^{10} and 2^{20} with items of length 128 bits, our PSU requires only 2.767 MB of communication. Compared with the state-of-the-art PSU proposed by Zhang et al. (USENIX Security 2023), there are $37\times$ shrink in communication and roughly $10 - 35\times$ speedup in the running time depending on the network environments.

Keywords: private set union; leveled fully homomorphic encryption; private equality test; standard security

1 Introduction

Private set union (PSU) is a cryptographic protocol that allows two parties, a sender and a receiver with respective input sets X and Y , to compute the union $X \cup Y$, without revealing anything else. It has numerous applications, such as cyber risk assessment [LV04, HLS⁺16, KRTW19], privacy-preserving data aggregation [BSMD10], and private ID [GMR⁺21] etc. However, most PSU protocols [KS05, Fri07, BA12, DC17, KRTW19, GMR⁺21] are designed for balanced inputs. These protocols typically perform only marginally better when one of the sets is much smaller than the other. Specifically, the communication cost of these protocols increases at least linearly with the size of the larger set.

In many real-world scenarios, the size of the sender’s input set is much smaller than that of the receiver, for instance, the sender is a mobile device with limited resources (e.g., battery, computing power, storage), while the receiver is a high-performance server. Additionally, the available bandwidth between the parties may be limited. As a concrete example, blacklist aggregation is a typical application of PSU [ZCL⁺23, JSZ⁺22, KRTW19], since the detection strategy may be deduced according to the IP addresses in the intersection. As reported in [RMY20], they have collected around 176 million blacklisted IP addresses over 23,483 autonomous systems, where the size of the IP blacklist varies greatly, with some blacklists listing more than 500,000 IP addresses and others listing fewer than 1,000 IP addresses. Therefore, aggregating IP blacklists with large size differences represents a representative use case of unbalanced PSU. However, most existing PSU protocols are not very efficient in dealing with the unbalanced case.

Over the last decade, there has been a significant amount of work on private set intersection (PSI) including both balanced case [PRTY19, PSTY19, CO18, FNO19, CCS18, HV17, RR17, PSZ18, PSZ14,

DCW13, HEK12, KMP⁺17] and unbalanced case [CLR17, KLS⁺17, CHLR18, RA18, CMdG⁺21], but little attention has been paid to PSU, particularly in the unbalanced setting. Recently, Jia et al. [JSZ⁺22] propose an unbalanced PSU*⁴ by using a shuffling technique, but their PSU* suffers from two primary drawbacks. Firstly, their PSU* does not satisfy standard security, since it *leaks the information of the intersection size* to the sender. Such information leakage could be critical for PSU. Consider an extreme case: the sender can get the intersection item if it inputs a singleton. Secondly, the communication complexity of their PSU* is *linear in the size of the larger set*. Chen et al. [CZZD22] shows how to tweak FHE-based PSI [CLR17] to an unbalanced PSU protocol. As noted by the authors, their PSU protocol only serves as a proof of concept since it reveals intersection size to the sender, and straightforward application of the optimization tricks due to [CLR17] will compromise the semi-honest security of the receiver. They left the standard security and efficient FHE-based PSU protocol in the unbalanced setting as a challenging problem.

Motivated by the above discussions, we ask the following question:

Is it possible to design a secure and fast unbalanced PSU protocol whose communication complexity is linear in the size of the smaller set, and logarithmic in the larger set?

1.1 Contributions

In this paper, we give an affirmative answer to the above question. Our contributions are summarized as follows:

- We present a generic construction of PSU protocol in the semi-honest model. The technical route is as follows:
 - We first propose a basic unbalanced PSU protocol based on the leveled fully homomorphic encryption (FHE). The basic protocol achieves communication linear in the small set and independent of the large set, but the deep depth of the homomorphic circuit leads to high computational costs.
 - Then, we use an array of optimization techniques following [CLR17, CHLR18, CMdG⁺21] to optimize the basic protocol, while the optimization might leak some information about the intersection. To remedy this issue, we introduce a new cryptographic protocol named permuted matrix Private Equality Test (pm-PEQT) to avoid information leakage, and propose two constructions of pm-PEQT. The first is based on Permute + Share [MS13, CGP20] and multi-point oblivious pseudorandom function (mp-OPRF). The second is based on the decisional Diffie-Hellman (DDH) assumption.
- By instantiating the generic construction, we obtain two secure and fast unbalanced PSU protocols whose communication complexity is *linear* in the size of the smaller set, and *logarithmic* in the larger set. Therefore, our protocols are particularly powerful when the set size of one party is much larger than that of the other.
- We implement our PSU protocols and compare them with the state-of-the-art PSU. Experiments show that our protocols are more efficient than all previous protocols in the unbalanced setting. For unbalanced sets size ($|X| = 2^{10}$, $|Y| = 2^{20}$) with 128-bit length items, our PSU protocol takes 2.767 MB of communication and about 18 seconds of computation with a single thread in LAN settings. There are roughly $37\times$ shrink in communication and $21\times$ speedup in the running time compared with the current most efficient PSU [ZCL⁺23]. In particular, the performance of our PSU protocols improves significantly in the case of low bandwidth. Our PSU requires 24.75 seconds in 1 Mbps bandwidth which is about $35\times$ faster than PSU of [ZCL⁺23].

1.2 Related Works

We give a brief survey of recent PSU protocols [KRTW19, GMR⁺21, JSZ⁺22, ZCL⁺23] as follows. For convenience, we define the set size $|X| = |Y| = n$ in the balanced case and $|X| = m \ll |Y| = n$ in the unbalanced case.

⁴ In this paper, we use PSU* to indicate a PSU protocol with information leakage.

Kolesnikov et al. [KRTW19] propose a PSU protocol based on the reverse private membership test (RPMT), in which the sender with input x interacts with the receiver holding a set Y , and the receiver can learn a bit indicating whether $x \in Y$, while the sender learns nothing. Their PSU is executed as follows. First, both parties run RPMT protocol, and then the receiver invokes OT protocol with the sender to obtain $\{x\} \cup Y$. For the balanced case, the protocol runs n RPMT instances independently and requires $O(n^2)$ communication and $O(n^2 \log^2 n)$ computation. By using the bucketing technique, two parties hash their sets into β bins and each bin consists of ρ items. A (n, n) -PSU⁵ is divided into β (ρ, ρ) -PSU. Consider the unbalanced case, the complexity is reduced to $O(n \log m)$ communication and $O(n \log m \log \log m)$ computation. However, as pointed in [JSZ⁺22], the bucketing technique leaks extra information to \mathcal{R} . More precisely, \mathcal{R} learns that some subsets (size ρ) hold the intersection items with high probability.

Garimella et al. [GMR⁺21] give a PSU protocol based on permuted characteristic functionality. We recall their PSU as follows. First, the sender holding a set X interacts with the receiver holding a set Y and invokes permuted characteristic protocol. As a result, the sender gets a random permutation π and the receiver obtains a vector $\mathbf{e} \in \{0, 1\}^n$, where $e_i = 1$ if $x_i \in Y$, else $e_i = 0$. After that, the receiver runs OT with the sender to obtain the union. Their protocol requires $O(n \log n)$ communication and $O(n \log n)$ computation in the balanced setting. Consider the unbalanced case, the complexity is reduced to $O(n + m \log m)$ communication and $O(n + m \log m)$ computation.

Jia et al. [JSZ⁺22] propose a PSU with the shuffling technique. Their protocol is executed as follows. First, the receiver hashes the set Y into Y_c by Cuckoo hashing and the sender hashes the set X by simple hashing. Then, two parties invoke the Permute + Share [MS13, CGP20] where the receiver shuffles Y_c by a permutation π chosen by the sender. As a result, the sender and the receiver obtain shuffled shares. After that, two parties invoke the mp-OPRF and the result is that the sender obtains the PRF key and the receiver obtains the PRF values. Then, the sender sends its PRF values to the receiver, and the receiver can test which items belong to the union. Finally, the receiver runs OT with the sender to get the union. Their PSU requires $O(n \log n)$ communication and $O(n \log n)$ computation. They also consider the unbalanced case and give an unbalanced PSU*, but their PSU* leaks the information of the intersection size to the sender and requires the communication complexity $O(n + m \log m)$ and $O(n + m \log m)$ computation.

Zhang et al. [ZCL⁺23] recently give a generic framework of PSU based on the multi-query reverse private membership test (mq-RPMT). We recall their PSU as follows. First, the sender holding a set X interacts with the receiver holding a set Y and runs mq-RPMT protocol. As a result, the sender obtains nothing and the receiver gets $\mathbf{b} \in \{0, 1\}^n$, satisfying that $b_i = 1$ if and only if $x_i \in Y$. Then, two parties run OT protocol to let the receiver get the union. They give two concrete mq-RPMT protocols. The first is based on the symmetric-key encryption and general 2PC. The second is based on the re-randomizable public-key encryption. Both constructions achieve $O(n)$ communication and $O(n)$ computation. Consider the unbalanced case, the complexity is $O(n + m)$ communication and $O(n + m)$ computation.

Table 1 provides a comparison of our protocols with the above mentioned PSU protocols in the unbalanced setting. n and m denote the size of the large set and the small set, respectively.

Protocols	Communication	Computation	Security
PSU* [KRTW19]	$O(n \log m)$	$O(n \log m \log \log m)$	Leaky
PSU [GMR ⁺ 21]	$O(n + m \log m)$	$O(n + m \log m)$	Standard
PSU [JSZ ⁺ 22]	$O(n + m \log m)$	$O(n + m \log m)$	Standard
PSU* [JSZ ⁺ 22]	$O(n + m \log m)$	$O(n + m \log m)$	Leaky
PSU [ZCL ⁺ 23]	$O(n + m)$	$O(n + m)$	Standard
Our PSU	$O(m \log n)$	$O(n + m \log n)$	Standard

Table 1: Comparisons of PSU in the semi-honest setting.

⁵ In this paper, we use (m, n) -PSU to indicate a PSU protocol where the sender's set size is m and the receiver's set size is n .

2 Technical Overview

We present an overview of our unbalanced PSU protocol as follows. First, we propose a basic PSU protocol based on leveled FHE, which is easy to understand but inefficient due to the deep depth of homomorphic circuits. Then, we try to improve the basic PSU by applying optimization techniques following [CLR17, CHLR18, CMdG+21] to reduce the depth of homomorphic circuits. However, the straightforward application of these optimization techniques leaks information about the intersection. To remedy this issue, we introduce a new cryptographic protocol called permuted matrix private equality test (pm-PEQT). Finally, we manage to give a generic construction of unbalanced PSU from leveled FHE and pm-PEQT. By instantiating the generic construction, we obtain secure and fast unbalanced PSU protocols in the semi-honest model. For convenience, we denote the parties in our PSU as the sender \mathcal{S} and the receiver \mathcal{R} , and their respective input sets as X and Y with $m = |X| \ll n = |Y|$.

2.1 Our Basic PSU Protocol

Our starting point is the basic FHE-based PSI protocol proposed by Chen et al. [CLR17], in which \mathcal{R} holding an item y interacts with \mathcal{S} holding a large set X , and \mathcal{R} can get the intersection $\{y\} \cap X$. We recall their construction as follows. First, \mathcal{R} uses its public key to encrypt y and sends the ciphertext $c \leftarrow \text{FHE.Enc}(y)$ to \mathcal{S} . Then, \mathcal{S} uses randomized product to hide its items by choosing random non-zero plaintext r and homomorphically computing $c' \leftarrow \text{FHE.Enc}(rf(y))$, where $f(x) = \prod_{x_i \in X} (x - x_i)$. After that, \mathcal{S} sends c' back and \mathcal{R} decrypts c' : if $rf(y) = 0$, it knows $y \in X$ and outputs $\{y\}$, else, it gets a random value and outputs \emptyset . The protocol achieves communication linear in the small set and independent of the large set, but it has high computational costs and deep homomorphic circuits since the degree of $f(x)$ is related to the large set size.

Basic PSU protocol. The functionality adjustment (PSI \rightarrow PSU) doesn't seem to be straightforward, since the randomized product $rf(y) = 0$ leaks the information of the intersection to the receiver. The challenge of solving above problem is to find a new randomization method that hides the information about the intersection and admits to checking which items belong to the union. We solve the problem by adding a random value r to randomize the polynomial value. In this way, the randomized value $r + f(y)$ leaks nothing to \mathcal{R} . If \mathcal{R} sends the result $r + f(y)$ to \mathcal{S} , \mathcal{S} can check whether the item y belongs to the union by verifying $r \stackrel{?}{=} r + f(y)$. Based on the requirements of the ideal functionality of PSU, which requires the receiver to output the union, we consider the opposite case of [CLR17], where \mathcal{R} holds the large set and \mathcal{S} holds the small set.

Suppose that \mathcal{S} has one item x and \mathcal{R} holding a large set Y gets the resulting union $\{x\} \cup Y$. We describe our basic PSU as follows. \mathcal{S} uses its public key to encrypt the item x and sends the ciphertext $c = \text{FHE.Enc}(x)$ to \mathcal{R} ; \mathcal{R} chooses random non-zero value r , and homomorphically computes $c' = \text{FHE.Enc}(r + f(x))$, where the polynomial $f(x) = \prod_{y_i \in Y} (x - y_i)$ and returns the new ciphertext to \mathcal{S} ; \mathcal{S} decrypts c' , gets the plaintext $r' = r + f(x)$ and returns r' back to \mathcal{R} ; \mathcal{R} sets $b = 0$ if $r' = r$, which means $x \in Y$, and $b = 1$ otherwise. Finally, \mathcal{R} invokes the OT protocol with \mathcal{S} to obtain the union $\{x\} \cup Y$. See Section 4 for details of our basic PSU protocol.

Similar to [CLR17], the basic PSU protocol achieves communication linear in the small set and independent of the large set, but the deep depth of the homomorphic circuit leads to high computational costs since the degree of $f(x)$ is related to the large set size.

2.2 Optimized PSU with Leakage

Chen et al. [CLR17] use an array of optimization techniques such as hashing, batching, windowing, partitioning, and modulus switching to improve their basic PSI protocol and obtain a fast PSI in the unbalanced setting. Therefore, we consider to adopt these optimization techniques to improve our basic PSU protocol as well.

First, we recall the optimized unbalanced PSI as follows: \mathcal{R} inserts the small set Y into Cuckoo hash table Y_c by Cuckoo hashing and each bin $Y_c[i]$ consists of one item. \mathcal{S} inserts the large set X into hash table X_b by

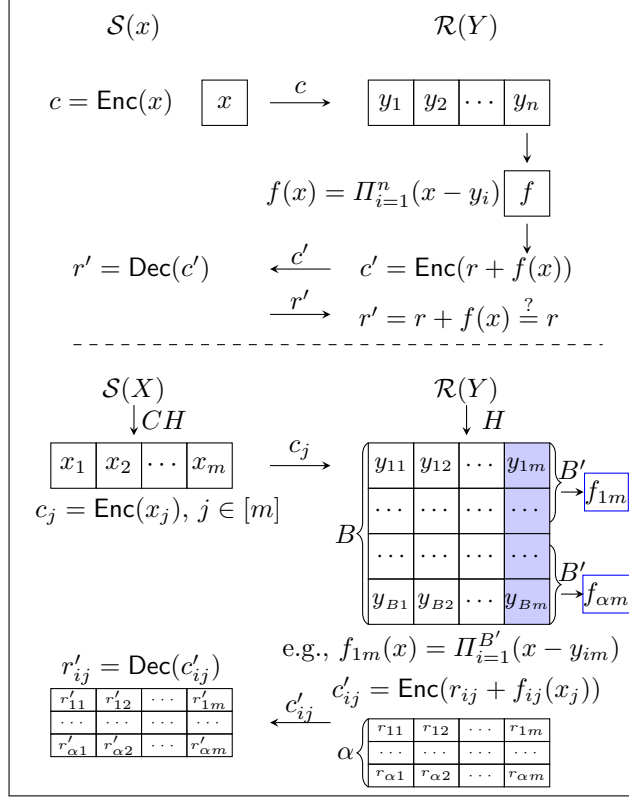


Fig. 1: The basic PSU (omit OT) and its optimizations. CH indicates Cuckoo Hashing and H indicates simple hashing.

simple hashing, where the i -th bin indicates $X_b[i]$ and each bin consists of B items. \mathcal{S} partitions each bin $X_b[i]$ into α subsets and each subset consists of $B' = B/\alpha$ items. Therefore, the (n, m) -PSI is divided into many $(B', 1)$ -PSI. For each $(B', 1)$ -PSI, \mathcal{S} encodes each subset (B' items) into a polynomial and randomizes it by multiplying a random value, then it homomorphically computes the new ciphertexts and sends them to \mathcal{R} . \mathcal{R} can decrypt the ciphertexts and obtain the intersection. Since the degree of the polynomial is related to the small subset size B' , each $(B', 1)$ -PSI has a low homomorphic circuit depth.

It is tempting to use the same optimization techniques [CLR17] to optimize our basic PSU as follows: \mathcal{S} inserts the small set X into a hash table X_c by Cuckoo hashing and each bin $X_c[i]$ consists of one item. \mathcal{R} inserts the large set Y into a hash table Y_b and each bin $Y_b[i]$ consists of B items. Then \mathcal{R} partitions each bin $Y_b[i]$ into α subsets and each subset consists of $B' = B/\alpha$ items. The (m, n) -PSU is divided into many $(1, B')$ -PSU. For each $(1, B')$ -PSU, \mathcal{R} encodes each subset (B' items) into a polynomial and randomizes it by adding a random value, then homomorphically computes the new ciphertexts and sends them to \mathcal{S} . After that, \mathcal{S} decrypts the ciphertexts and sends the plaintexts back. Finally, \mathcal{R} checks which items belong to the set union, and invokes OT with \mathcal{S} to get them. We show our basic PSU (omit OT) and its optimizations in Figure 1.

However, unlike PSI [CLR17], the optimization techniques for PSI are not suitable for our PSU. This is because a large PSI can be divided into many small PSI, and the receiver can combine all small set intersections into the output securely, but if we divide a (m, n) -PSU into many $(1, B')$ -PSU directly, this causes information leakage about the intersection. We show the comparison of PSI and our optimized PSU (omit OT) with leakage in Figure 2. Note that in the (m, n) -PSU, from the view of \mathcal{R} , any item in the set Y could be an item in $X \cap Y$. However, in the above optimized PSU*, \mathcal{R} learns some subsets with size B' have the item in $X \cap Y$. Moreover, if \mathcal{S} returns its decrypted results r' to \mathcal{R} directly. \mathcal{R} can check which

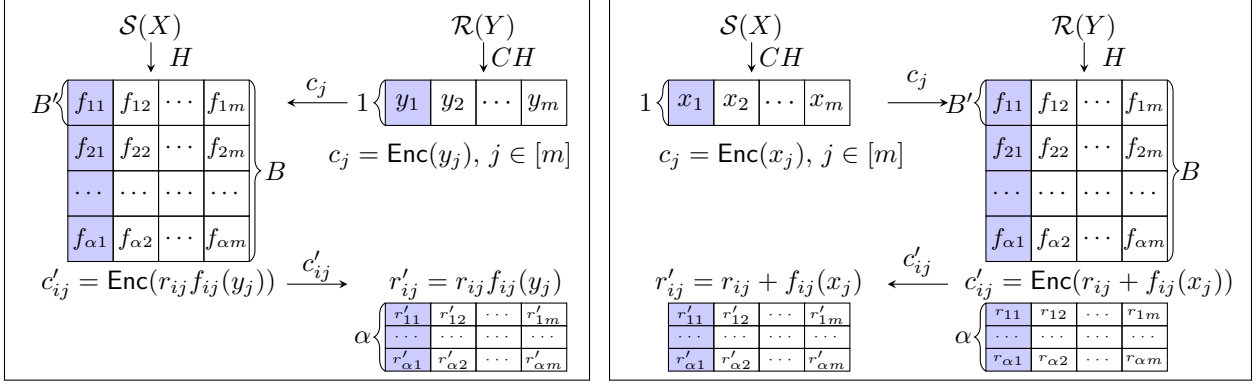


Fig. 2: Comparison of PSI [CLR17] (left) and optimized PSU (omit OT) with leakage (right)

items of X belong to the set union. This also leaks the information of $X \cap Y$. Because there are α subsets with size B' in one bin if $f(x) = 0$ in one subset, \mathcal{R} gets $f'(x) \neq 0$ in other subsets, which causes \mathcal{R} could compute the intersection items with sufficient polynomial values. For example, in Figure 2 (right), in the first column, if $r_{11} = r'_{11}$, this means $x_1 \in Y^6$ and $x_1 \in \{y_{11}, \dots, y_{B'1}\}$, but $x_1 \notin \{y_{(B'+1)1}, \dots, y_{B1}\}$. \mathcal{R} gets the rest nonzero polynomial values $f_{21}(x_1), \dots, f_{\alpha 1}(x_1)$ and it could compute x_1 from them.

Based on the above analysis, the main challenge is how to optimize our basic PSU without causing information leakage. More precisely, for \mathcal{S} holds a matrix $\mathbf{R}' = [r'_{ij}]$ and \mathcal{R} holds a matrix $\mathbf{R} = [r_{ij}]$, $i \in [\alpha], j \in [m]$, we need to overcome the following two difficulties:

- The receiver is able to check $r_{ij} \stackrel{?}{=} r'_{ij}$ for all $i \in [\alpha], j \in [m]$ without knowing the value r'_{ij} held by the sender.
- The receiver is able to check $r_{ij} \stackrel{?}{=} r'_{ij}$ for all $i \in [\alpha], j \in [m]$ without knowing the positions i, j of r_{ij} .

2.3 Permuted Matrix Private Equality Test

In order to address the above difficulties, we introduce a new cryptographic protocol named permuted matrix private equality test (pm-PEQT) which can be seen as an extension of private equality test (PEQT). In the PEQT, a receiver who has an input string x interacts with a sender holding an input string y , and the result is that the receiver learns a bit indicating whether $x = y$ merely, whereas the sender learns nothing. In our pm-PEQT, a sender holding a matrix $\mathbf{R}'_{\alpha \times m}$ and a matrix permutation $\pi = (\pi_c, \pi_r)$ interacts with a receiver holding a matrix $\mathbf{R}_{\alpha \times m}$. As a result, the receiver learns (only) the bit matrix $\mathbf{B}_{\alpha \times m}$ indicating that $b_{ij} = 1$ if $r_{\pi(ij)} = r'_{\pi(ij)}$ and $b_{ij} = 0$ otherwise, for $i \in [\alpha], j \in [m]$, while the sender learns nothing about \mathbf{R} . Compared with PEQT, pm-PEQT admits a matrix private equality test with *positions permutation*.

Constructions of pm-PEQT. pm-PEQT can not be easily built from PEQT by running many PEQT instances in parallel. The difficulty is to shuffle the receiver's items without knowing the permutation of the sender. We give two constructions of pm-PEQT as follows.

The first construction is based on Permute + Share [MS13, CGP20] and mp-OPRF [CM20]. Informally, \mathcal{S} and \mathcal{R} invoke the ideal Permute + Share functionality \mathcal{F}_{PS} twice for the columns and rows of \mathbf{R} , respectively: \mathcal{R} inputs the matrix $\mathbf{R}_{\alpha \times m}$ and \mathcal{S} inputs the permutation $\pi = (\pi_c, \pi_r)$, where π_c (over $[m]$) and π_r (over $[\alpha]$) are two sub-permutations for columns and rows. As a result, \mathcal{R} obtains the shuffled matrix shares \mathbf{S}_π and \mathcal{S} obtains the shuffled matrix shares \mathbf{S}'_π , where $s_{\pi(ij)} \oplus s'_{\pi(ij)} = r_{\pi(ij)}$, $i \in [\alpha], j \in [m]$. After that, both parties

⁶ In j -th column, as long as there is a position i , such that $r_{ij} = r'_{ij}$, $x_j \in Y$. Meanwhile, at most one position is equal in each column.

invoke mp-OPRF functionality $\mathcal{F}_{\text{mp-OPRF}}$: \mathcal{R} inputs shuffled shares \mathbf{S}_π and obtains $F_k(s_{\pi(ij)}), i \in [\alpha], j \in [m]$, and \mathcal{S} gets the PRF key k . Furthermore, \mathcal{S} permutes the matrix \mathbf{R}' by $\pi = (\pi_c, \pi_r)$ and gets $\mathbf{R}'_\pi = [r'_{\pi(ij)}]$, then \mathcal{S} computes all PRF values $F_k(r'_{\pi(ij)} \oplus s'_{\pi(ij)}), i \in [\alpha], j \in [m]$ and sends them to \mathcal{R} . Finally, \mathcal{R} sets $b_{ij} = 1$, if $F_k(s_{\pi(ij)}) = F_k(r'_{\pi(ij)} \oplus s'_{\pi(ij)})$, else, sets $b_{ij} = 0$. \mathcal{R} defines the bit matrix $\mathbf{B} = [b_{ij}], i \in [\alpha], j \in [m]$. The details of this construction are deferred to Section 5.1. The communication complexity of our pm-PEQT based on Permute + Share and mp-OPRF is $O(\alpha m \log \alpha m)$.

The second construction is based on the DDH assumption. Let \mathbb{G} be a cyclic group with order q . First, \mathcal{R} and \mathcal{S} choose random values $a, b \leftarrow \mathbb{Z}_q$ and compute $v_{ij} = H(r_{ij})^a, v'_{ij} = H(r'_{ij})^b$ for all $i \in [\alpha], j \in [m]$, respectively, where $H(\cdot)$ is modeled as a random oracle and the output is a group element in \mathbb{G} . Let $\mathbf{V} = [v_{ij}]$ and $\mathbf{V}' = [v'_{ij}], i \in [\alpha], j \in [m]$. \mathcal{R} sends \mathbf{V} to \mathcal{S} . Then \mathcal{S} computes $v''_{ij} = (v_{ij})^b$ and defines $\mathbf{V}'' = [v''_{ij}], i \in [\alpha], j \in [m]$. \mathcal{S} shuffles \mathbf{V}'' and \mathbf{V}' by same permutation $\pi = (\pi_c, \pi_r)$ and gets $\mathbf{V}''_\pi = \pi(\mathbf{V}'')$, $\mathbf{V}'_\pi = \pi(\mathbf{V}')$, where $v''_{\pi(ij)} = \pi(v''_{ij}), v'_{\pi(ij)} = \pi(v'_{ij})$. \mathcal{S} sends permuted matrices \mathbf{V}''_π and \mathbf{V}'_π to \mathcal{R} . Finally, for i -th row and j -column in \mathbf{V}''_π and \mathbf{V}'_π , if $v''_{\pi(ij)} = v'_{\pi(ij)}$, \mathcal{R} defines $b_{ij} = 1$ and $b_{ij} = 0$ otherwise. Hence, \mathcal{R} gets a bit matrix $\mathbf{B} = [b_{ij}], i \in [\alpha], j \in [m]$. The details of this construction are deferred to section 5.2. The communication complexity of our DDH-based pm-PEQT is $O(\alpha m)$.

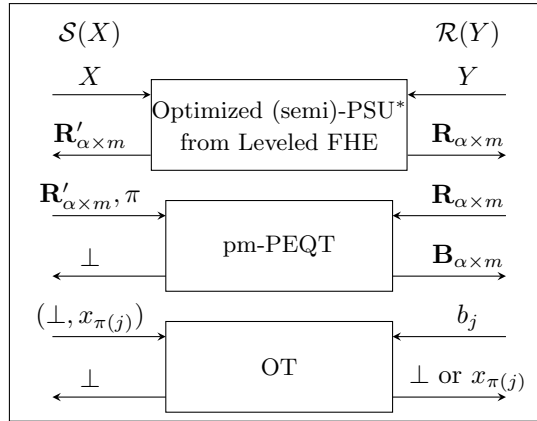


Fig. 3: Core design idea of our full PSU protocol

2.4 Our Full PSU Protocol

Now, we are ready to describe our full PSU protocol. We provide the high-level technical overview for our generic construction of PSU in Figure 3 and the details are as follows.

First, we construct a semi-finished FHE-based optimized PSU* in which \mathcal{S} does not send its decrypted results $\mathbf{R}'_{\alpha \times m}$ to \mathcal{R} . In this step, \mathcal{S} holding a small set X interacts with \mathcal{R} holding a large set Y . The result is that \mathcal{R} outputs a matrix $\mathbf{R}_{\alpha \times m} = [r_{ij}]$, and \mathcal{S} outputs a matrix $\mathbf{R}'_{\alpha \times m} = [r'_{ij}]$, where α denotes the number of partitions, m denotes the number of bins, r_{ij} denotes the random values used to hide each polynomial and r'_{ij} denotes the decrypted results, $i \in [\alpha], j \in [m]$. Note that for all $i \in [\alpha]$ in same j -th column, if all $r'_{ij} \neq r_{ij}$, we have $x_j \notin Y$, else, $x_j \in Y$.

Then, two parties invoke the pm-PEQT functionality: \mathcal{S} inputs the matrix \mathbf{R}' and a permutation $\pi = (\pi_c, \pi_r)$ ⁷ and \mathcal{R} inputs the matrix \mathbf{R} . As a result, \mathcal{R} gets a bit matrix $\mathbf{B} = [b_{ij}]$, where $b_{ij} = 1$, if and only if $r_{\pi(ij)} = r'_{\pi(ij)}$ for $i \in [\alpha], j \in [m]$. \mathcal{R} computes a bit vector $\mathbf{b} = [b_j], j \in [m]$, for all $i \in [\alpha]$, if $b_{ij} = 0$, sets

⁷ Each column of the matrix $\mathbf{R}_{\alpha \times m}$ corresponds to the same item, so the permutation for the matrix requires that the columns are consistent.

$b_j = 1$, else, sets $b_j = 0$. \mathcal{S} permutes the Cuckoo hash table X_c by π_c and gets $\pi_c(X_c) = [x_{\pi_c(1)}, \dots, x_{\pi_c(m)}]$. We note that if $b_j = 1$, we have $x_{\pi_c(j)} \notin Y$, else $x_{\pi_c(i)} \in Y$, $j \in [m]$.

Finally, for $j \in [m]$, two parties execute OT protocol as follows. In the j -th OT instance, \mathcal{S} inputs $(\perp, x_{\pi_c(j)})$, and \mathcal{R} inputs b_j , $j \in [m]$. As a result, \mathcal{R} obtains all $x_{\pi_c(j)}$ for $b_j = 1$. \mathcal{R} outputs the union $Y \cup \{x_{\pi_c(j)}\}_{j \in [m]}$.

3 Preliminaries

3.1 Notation

For $n \in \mathbb{N}$, let $[n]$ denote the set $\{1, 2, \dots, n\}$. 1^λ denotes the string of λ ones. We use κ and λ to indicate the computational and statistical security parameters, respectively. If S is a set, $s \leftarrow S$ indicates sampling s from S at random. We denote vectors by lower-case bold letters, e.g., \mathbf{s} . We denote matrices by upper-case bold letters, e.g., \mathbf{S} . We write $\mathbf{S} = [s_{ij}]$ to denote the matrix \mathbf{S} with i -th row and j -th column element s_{ij} . We write $\mathbf{S}_{n \times m}$ to denote the matrix \mathbf{S} has n rows and m columns. For a permutation π over n items, we write $\{s_{\pi(1)}, \dots, s_{\pi(n)}\}$ to denote $\pi(\{s_1, \dots, s_n\})$, where $s_{\pi(i)}$ indicates the i -th element after the permutation. For a column permutation π_c (or, row permutation π_r) on a matrix $\mathbf{S} = [s_{ij}]$, we write \mathbf{S}_{π_c} (or, \mathbf{S}_{π_r}) to denote that $\pi_c(\mathbf{S}) = [s_{\pi_c(ij)}]$ (or, $\pi_r(\mathbf{S}) = [s_{\pi_r(ij)}]$) is the permuted matrix, where $s_{\pi_c(ij)}$ (or, $s_{\pi_r(ij)}$) indicates the i -th row and j -th column element after the permutation.

3.2 Private Set Union

PSU is a special case of secure two-party computation. We describe the ideal functionality of PSU in Figure 4.

Parameters: Set sizes m and n .

Functionality:

1. Wait for an input $X = \{x_1, \dots, x_m\} \subseteq \{0, 1\}^*$ from the sender, and an input $Y = \{y_1, \dots, y_n\} \subseteq \{0, 1\}^*$ from the receiver.
2. Give output $X \cup Y$ to the receiver.

Fig. 4: Ideal functionality $\mathcal{F}_{\text{PSU}}^{m,n}$ for private set union

3.3 Building Blocks

We briefly review the main technical tools including simple hashing, Cuckoo hashing, leveled fully homomorphic encryption, oblivious transfer, multi-point oblivious PRF, and Permute + Share.

Simple hashing. In the simple hashing [PSZ14], the hash table consists of m bins B_1, \dots, B_m . Hashing is done by mapping each element x to a bin $B_{h(x)}$ using a hash function $h : \{0, 1\}^* \rightarrow [m]$ that was chosen uniformly at random and independently of the input elements. According to the following inequality [MR95], the maximum bin size B can be set to ensure that no bin will contain more than B items except with probability $2^{-\lambda}$ when hashing n items into m bins.

$$\Pr[\exists \text{ bin size} \geq B] \leq m \left[\sum_{i=B}^n \binom{n}{i} \cdot \left(\frac{1}{m}\right)^i \cdot \left(1 - \frac{1}{m}\right)^{n-i} \right]$$

Cuckoo hashing. Cuckoo hashing can be used to build dense hash tables by many hash functions [PR01, DM03, FPSS03, PSZ18]. Following [CLR17, CHLR18], we use three hash functions and adjust the number of items and table size to reduce the stash size to 0 while achieving a hashing failure probability of $2^{-\lambda}$.

Leveled fully homomorphic encryption. The leveled fully homomorphic encryption supports circuits of a certain bounded depth. Following [CLR17], our protocols require that the leveled FHE satisfies IND-CPA secure with circuit privacy [BPMW16]. We use an array of optimization techniques of FHE as [CLR17, CHLR18, CMdG⁺21], such as batching, windowing, partitioning, and modulus switching to significantly reduce the depth of the homomorphic circuit. We review these optimizations in Appendix A. For the implementation, we use the homomorphic encryption library SEAL which implements the BFV scheme [FV12] following [CLR17, CHLR18, CMdG⁺21].

Oblivious transfer. Oblivious transfer (OT) [Rab05] is a central cryptographic primitive in the area of secure computation. In the 1-out-of-2 OT, a sender with two input strings (x_0, x_1) interacts with a receiver who has an input choice bit b . The result is that the receiver learns x_b without learning anything about x_{1-b} , while the sender learns nothing about b . Ishai et al. [IKNP03] introduced the OT extension that allows for a large number of OT executions at the cost of computing a small number of public-key operations. We recall the 1-out-of-2 oblivious transfer functionality \mathcal{F}_{OT} in Figure 5.

Parameters: Two parties: \mathcal{S} and \mathcal{R} .
Functionality \mathcal{F}_{OT} :

1. Wait for input $\{x_0, x_1\}$ from \mathcal{S} . Wait for input $b \in \{0, 1\}$ from \mathcal{R} .
2. Give x_b to \mathcal{R} .

Fig. 5: 1-out-of-2 oblivious transfer functionality

Multi-point oblivious pseudorandom function. An oblivious pseudorandom function (OPRF) allows a receiver to input x and learns the PRF value $F_k(x)$, and the key k is known to a sender. Pinkas et al. [PRTY19] propose multi-point OPRF (mp-OPRF) and realize efficient PSI protocols. Recently, Chase and Miao [CM20] propose a more efficient mp-OPRF based on oblivious transfer extension. In the mp-OPRF, the receiver inputs $\{x_1, x_2, \dots, x_n\}$ and learns all PRF values $\{F_k(x_1), F_k(x_2), \dots, F_k(x_n)\}$, and the sender gets the key k . We recall the mp-OPRF functionality $\mathcal{F}_{\text{mp-OPRF}}$ in Figure 6.

Parameters: A PRF F . Two parties: \mathcal{S} and \mathcal{R} .
Functionality $\mathcal{F}_{\text{mp-OPRF}}$:

1. Wait for input $\{x_1, \dots, x_n\}$ from \mathcal{R} .
2. Sample a random PRF key k and give it to \mathcal{S} . Give $\{F_k(x_1), \dots, F_k(x_n)\}$ to \mathcal{R} .

Fig. 6: mp-OPRF functionality

Permute + Share. We recall the Permute + Share functionality \mathcal{F}_{PS} [MS13, CGP20] in Figure 7. Roughly speaking, in the Permute + Share protocol, P_0 inputs a set $X = \{x_1, \dots, x_n\}$ of size n and P_1 chooses a permutation π on n items. The result is that P_0 learns the shuffled shares $\{s_{\pi(1)}, \dots, s_{\pi(n)}\}$ and P_1 learns the other shuffled shares $\{s'_{\pi(1)}, \dots, s'_{\pi(n)}\}$, where $x_{\pi(i)} = s_{\pi(i)} \oplus s'_{\pi(i)}$, $i \in [n]$.

Parameters: Two parties: P_0 and P_1 . Set size n for P_0 .

Functionality \mathcal{F}_{PS} :

1. Wait for input $X = \{x_1, \dots, x_n\}$ from P_0 , abort if $|X| \neq n$. Wait for input a permutation π from P_1 , abort if π is not a permutation on n items.
2. Give output shuffled shares $\{s_{\pi(1)}, \dots, s_{\pi(n)}\}$ to P_0 , and another shuffled shares $\{s'_{\pi(1)}, \dots, s'_{\pi(n)}\}$ to P_1 , where $x_{\pi(i)} = s_{\pi(i)} \oplus s'_{\pi(i)}$, $i \in [n]$.

Fig. 7: Permute + Share functionality

4 The Basic PSU Protocol

We describe our basic PSU protocol in Figure 8 as a strawman protocol.

Input: The sender inputs set X of size $m = |X|$ and the receiver inputs set Y of size $n = |Y|$, where m and n are public.

Output: The receiver outputs $X \cup Y$. The sender outputs \perp .

- (1) **[Setup]** \mathcal{S} generates a public-secret key pair and keeps the secret key itself.
- (2) **[Set encryption]** \mathcal{S} encrypts each item $x_i \in X$, $c_i = \text{FHE.Enc}(x_i)$, $i \in [m]$ and sends (c_1, \dots, c_m) to \mathcal{R} .
- (3) **[Computation]** For each c_i , \mathcal{R}
 - (a) samples a random non-zero value r_i ;
 - (b) homomorphically computes $c'_i = \text{FHE.Enc}(f(x_i) + r_i)$, where $f(x) = \prod_{y_i \in Y} (x - y_i)$.
 - (c) sends c'_i , $i \in [m]$ to \mathcal{S} .
- (4) **[Decryption]** \mathcal{S} decrypts c'_i , $i \in [m]$ to $r'_i = f(x_i) + r_i$ and sends them to \mathcal{R} .
- (5) **[Output]** \mathcal{R} checks all plaintexts and defines a bit vector $\mathbf{b} = [b_i]$, $i \in [m]$. If $r'_i = r_i$, \mathcal{R} sets $b_i = 0$, otherwise, sets $b_i = 1$. Then, both parties invoke OT protocol, where \mathcal{R} inputs the bit b_i and \mathcal{S} inputs (\perp, x_i) , $i \in [m]$. As a result, \mathcal{R} obtains x_i corresponding to $b_i = 1$. Finally, \mathcal{R} outputs $X \cup Y := Y \cup \{x_i | b_i = 1\}_{i \in [m]}$.

Fig. 8: Basic PSU protocol

We prove its semi-honest security in the following theorem.

Theorem 1. *The PSU protocol described in Figure 8 is secure in the \mathcal{F}_{OT} -hybrid model, in the presence of semi-honest adversaries, provided that the fully homomorphic encryption scheme is IND-CPA secure with circuit privacy.*

Proof. We construct $\text{Sim}_{\mathcal{S}}$ and $\text{Sim}_{\mathcal{R}}$ to simulate the views of corrupt \mathcal{S} and corrupt \mathcal{R} respectively and argue the indistinguishability of the produced transcript from the real execution.

Corrupt Sender. $\text{Sim}_{\mathcal{S}}(X)$ simulates the view of corrupt \mathcal{S} as follows: it encrypts m random values. Then, it invokes $\text{Sim}_{\text{OT}}^{\mathcal{S}}(\perp, x_i)$, $i \in [m]$ 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. The plaintexts are randomized in the real view which is indistinguishable from the random values in the simulated view. The FHE satisfies the circuit privacy which hides the computational circuit in step 3. The view produced by the underlying OT simulator is indistinguishable from the real view. Thus, the simulation is indistinguishable from the real view.

Corrupt Receiver. $\text{Sim}_{\mathcal{R}}(Y, X \cup Y)$ simulates the view of corrupt \mathcal{R} as follows: it simulates the ciphertexts by encrypting m random values. $\text{Sim}_{\mathcal{R}}$ sets $\hat{X} = (X \cup Y) \setminus Y$ and pads \hat{X} with \perp into m items and permutes all items randomly. It computes the polynomial $f(y) = \prod_{y_i \in Y} (y - y_i)$ and the random values $\mathbf{r} = [r_i]$, $i \in [m]$ used to randomize the polynomial. Then, for $\hat{x}_i \neq \perp$, $\text{Sim}_{\mathcal{R}}$ defines $r'_i := f(\hat{x}_i) + r_i$, else, it defines $r'_i := r_i$, and appends $\mathbf{r}' = [r'_i]$, $i \in [m]$ to the view. If $\hat{x}_i = \perp$, it sets $b_i = 0$, else, $b_i = 1$. Then $\text{Sim}_{\mathcal{R}}$ invokes $\text{Sim}_{\text{OT}}^{\mathcal{R}}(b_i, \hat{x}_i)$ for $i \in [m]$ and appends the output to the view.

We argue that the outputs of $\text{Sim}_{\mathcal{R}}$ are indistinguishable from the real view of \mathcal{R} by the following hybrids:
Hyb₀: \mathcal{R} 's view in the real protocol.

Hyb₁: Same as **Hyb₀** except that the ciphertexts in the step 2 are replaced by encrypting m random values generated by $\text{Sim}_{\mathcal{R}}$. Since the fully homomorphic encryption scheme is IND-CPA secure, the above simulation is indistinguishable from the real view.

Hyb₂: Same as **Hyb₁** except that $\text{Sim}_{\mathcal{R}}$ runs the \mathcal{F}_{OT} simulator to produce the simulated view for \mathcal{R} . The security of OT protocol guarantees the view in simulation is computationally indistinguishable from the view in the real protocol. The hybrid is the view output by $\text{Sim}_{\mathcal{R}}$.

5 Permuted Matrix Private Equality Test

In this section, we introduce the ideal functionality of pm-PEQT in Figure 9, and then give two efficient constructions of pm-PEQT in the semi-honest model.

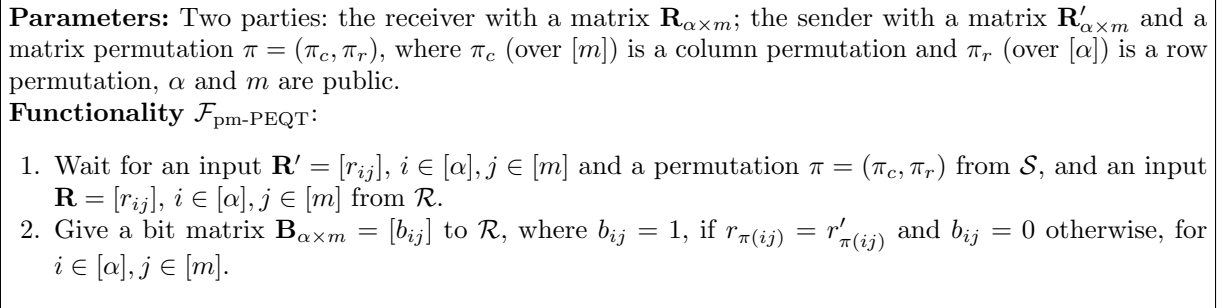


Fig. 9: Permuted matrix private equality test

5.1 pm-PEQT from Permute + Share and mp-OPRF

The first construction of pm-PEQT is based on the Permute + Share [CGP20] and mp-OPRF [CM20] as described in Figure 10. For notational convenience, we use $\text{pm-PEQT}_{\text{sym}}$ to denote pm-PEQT based on the Permute + Share and mp-OPRF, since the main computational costs of this protocols are symmetric-key operations.

Input: The receiver inputs a matrix $\mathbf{R}_{\alpha \times m}$. The sender inputs a matrix $\mathbf{R}'_{\alpha \times m}$, and a permutation $\pi = (\pi_c, \pi_r)$, where π_c (over $[m]$) and π_r (over $[\alpha]$) are two sub-permutations for columns and rows.

Output: The receiver outputs a matrix \mathbf{B} . The sender outputs \perp .

1. \mathcal{S} and \mathcal{R} invoke the ideal Permute + Share functionality \mathcal{F}_{PS} twice.
 - (a) Both parties permute and share the columns of \mathbf{R} , where each column of \mathbf{R} can be seen as an item. \mathcal{R} inputs each column of \mathbf{R} and \mathcal{S} inputs the permutation π_c . As a result, \mathcal{R} gets $\mathbf{S}_{\pi_c} = [s_{\pi_c(ij)}]$ and \mathcal{S} gets $\mathbf{S}'_{\pi_c} = [s'_{\pi_c(ij)}]$, where $s_{\pi_c(ij)} \oplus s'_{\pi_c(ij)} = r_{\pi_c(ij)}$.
 - (b) Both parties permute and share the rows of \mathbf{S}_{π_c} , where each row of \mathbf{S}_{π_c} can be seen as an item. \mathcal{R} inputs each row of \mathbf{S}_{π_c} and \mathcal{S} inputs the permutation π_r . As a result, \mathcal{R} gets $\mathbf{S}_{\pi_r} = [s_{\pi_r(ij)}]$ and \mathcal{S} gets $\mathbf{S}'_{\pi_r} = [s'_{\pi_r(ij)}]$, where $s_{\pi_r(ij)} \oplus s'_{\pi_r(ij)} = s_{\pi_c(ij)}$.
 - (c) \mathcal{R} defines the shuffled matrix shares $\mathbf{S}_\pi = \mathbf{S}_{\pi_r}$ and \mathcal{S} defines the shuffled matrix shares $\mathbf{S}'_\pi = \pi_r(\mathbf{S}'_{\pi_c}) \oplus \mathbf{S}'_{\pi_r}$, where $s_{\pi(ij)} \oplus s'_{\pi(ij)} = r_{\pi(ij)}$, $i \in [\alpha], j \in [m]$.
2. Both parties invoke mp-OPRF functionality $\mathcal{F}_{\text{mp-OPRF}}$. \mathcal{R} inputs shuffled shares \mathbf{S}_π and obtains the outputs $F_k(s_{\pi(ij)})$, $i \in [\alpha], j \in [m]$. \mathcal{S} obtains the key k .
3. \mathcal{S} computes $F_k(r'_{\pi(ij)} \oplus s'_{\pi(ij)})$, $i \in [\alpha], j \in [m]$ and sends them to \mathcal{R} .
4. \mathcal{R} sets $b_{ij} = 1$, if $F_k(s_{\pi(ij)}) = F_k(r'_{\pi(ij)} \oplus s'_{\pi(ij)})$, else, $b_{ij} = 0$, and gets a bit matrix $\mathbf{B} = [b_{ij}]$, $i \in [\alpha], j \in [m]$.

Fig. 10: pm-PEQT from Permute + Share and mp-OPRF

Theorem 2. *The construction of Figure 10 securely implements functionality $\mathcal{F}_{\text{pm-PEQT}}$ in the $(\mathcal{F}_{\text{PS}}, \mathcal{F}_{\text{mp-OPRF}})$ -hybrid model, in the presence of semi-honest adversaries.*

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. $\text{Sim}_{\mathcal{S}}(\mathbf{R}', \pi = (\pi_c, \pi_r))$ simulates the view of corrupt \mathcal{S} as follows: $\text{Sim}_{\mathcal{S}}$ randomly chooses \mathbf{S}'_{π_c} and invokes $\text{Sim}_{\text{PS}}^{\mathcal{S}}(\pi_c, \mathbf{S}'_{\pi_c})$ and appends the output to the view. $\text{Sim}_{\mathcal{S}}$ randomly chooses \mathbf{S}'_{π_r} and invokes $\text{Sim}_{\text{PS}}^{\mathcal{S}}(\pi_r, \mathbf{S}'_{\pi_r})$ and appends the output to the view. Then, $\text{Sim}_{\mathcal{S}}$ randomly selects a key k of PRF and invokes $\text{Sim}_{\text{mp-OPRF}}^{\mathcal{S}}(k)$ and appends the output to the view.

We argue that the outputs of $\text{Sim}_{\mathcal{S}}$ are indistinguishable from the real view of \mathcal{S} by the following hybrids:

Hyb₀: \mathcal{S} 's view in the real protocol.

Hyb₁: Same as **Hyb₀** except that the output of \mathcal{F}_{PS} is replaced by $\mathbf{S}'_{\pi_c}, \mathbf{S}'_{\pi_r}$ chosen by $\text{Sim}_{\mathcal{S}}$, and $\text{Sim}_{\mathcal{S}}$ runs the \mathcal{F}_{PS} simulator to produce the simulated view for \mathcal{S} . The security of Permute + Share guarantees the view in simulation is computationally indistinguishable from the view in the real protocol.

Hyb₂: Same as **Hyb₁** except that the output key of $\mathcal{F}_{\text{mp-OPRF}}$ is replaced by the k chosen by $\text{Sim}_{\mathcal{S}}$, and $\text{Sim}_{\mathcal{S}}$ runs the $\mathcal{F}_{\text{mp-OPRF}}$ simulator to produce the simulated view for \mathcal{S} . The security of mp-OPRF guarantees the view in simulation is computationally indistinguishable from the view in the real protocol. The hybrid is the view output by $\text{Sim}_{\mathcal{S}}$.

Corrupt Receiver. $\text{Sim}_{\mathcal{R}}(\mathbf{R}, \mathbf{B} = [b_{ij}])$ simulates the view of corrupt \mathcal{R} as follows: $\text{Sim}_{\mathcal{R}}$ chooses \mathbf{S}_{π_c} and invokes the simulator $\text{Sim}_{\text{PS}}^{\mathcal{R}}(\mathbf{R}, \mathbf{S}_{\pi_c})$ and appends the output to the view. $\text{Sim}_{\mathcal{R}}$ chooses \mathbf{S}_{π_r} and invokes $\text{Sim}_{\text{PS}}^{\mathcal{R}}(\mathbf{S}_{\pi_c}, \mathbf{S}_{\pi_r})$ and appends the output to the view. $\text{Sim}_{\mathcal{R}}$ randomly selects u_{ij} , $i \in [\alpha], j \in [m]$ and invokes $\text{Sim}_{\text{mp-OPRF}}^{\mathcal{R}}(u_{ij})$ and appends the output to the view. Finally, for all $i \in [\alpha], j \in [m]$, $\text{Sim}_{\mathcal{R}}$ sets $v_{ij} = u_{ij}$ if $b_{ij} = 1$, else, it chooses v_{ij} randomly and appends all v_{ij} to the view.

The view generated by $\text{Sim}_{\mathcal{R}}$ is indistinguishable from a real view of \mathcal{R} by the following hybrids:

Hyb₀: \mathcal{R} 's view in the real protocol.

Hyb₁: Same as **Hyb₀** except that the output of \mathcal{F}_{PS} is replaced by $\mathbf{S}_{\pi_c}, \mathbf{S}_{\pi_r}$ chosen by $\text{Sim}_{\mathcal{R}}$, and $\text{Sim}_{\mathcal{R}}$ runs the \mathcal{F}_{PS} simulator to produce the simulated view for \mathcal{R} . The security of Permute + Share guarantees the view in simulation is computationally indistinguishable from the view in the real protocol.

Hyb₂: Same as **Hyb₁** except that the output PRF values of $\mathcal{F}_{\text{mp-OPRF}}$ is replaced by u_{ij} , $i \in [\alpha]$, $j \in [m]$, and all PRF values in the last step is replaced by the v_{ij} , chosen by $\text{Sim}_{\mathcal{R}}$ randomly, and $\text{Sim}_{\mathcal{R}}$ runs the $\mathcal{F}_{\text{mp-OPRF}}$ simulator to produce the simulated view for \mathcal{R} . The security of mp-OPRF guarantees the view in simulation is computationally indistinguishable from the view in the real protocol.

Input: The receiver inputs a matrix $\mathbf{R}_{\alpha \times m}$. The sender inputs a matrix $\mathbf{R}'_{\alpha \times m}$, and a permutation $\pi = (\pi_c, \pi_r)$ where π_c (over $[m]$) and π_r (over $[\alpha]$) are two sub-permutation for columns and rows. \mathbb{G} is a cyclic group with order q .

Output: The receiver outputs a matrix \mathbf{B} . The sender outputs \perp .

1. \mathcal{R} chooses a random value $a \leftarrow \mathbb{Z}_q$ and computes $v_{ij} = H(r_{ij})^a$ for all $i \in [\alpha]$, $j \in [m]$, where $H(\cdot)$ is modeled as a random oracle and the output is a group element in \mathbb{G} . Let $\mathbf{V} = [v_{ij}]$, $i \in [\alpha]$, $j \in [m]$. \mathcal{R} sends \mathbf{V} to \mathcal{S} .
2. \mathcal{S} chooses a random value $b \leftarrow \mathbb{Z}_q$ and computes $v'_{ij} = H(r'_{ij})^b$ for all $i \in [\alpha]$, $j \in [m]$. Let $\mathbf{V}' = [v'_{ij}]$, $i \in [\alpha]$, $j \in [m]$. Then, \mathcal{S} computes $v''_{ij} = (v_{ij})^b$ and lets $\mathbf{V}'' = [v''_{ij}]$. After that, \mathcal{S} shuffles \mathbf{V}'' and \mathbf{V}' by same permutation $\pi = (\pi_c, \pi_r)$ and gets $\mathbf{V}''_{\pi} = \pi(\mathbf{V}'')$, $\mathbf{V}'_{\pi} = \pi(\mathbf{V}')$, where $v''_{\pi(ij)} = \pi(v''_{ij})$, $v'_{\pi(ij)} = \pi(v'_{ij})$. \mathcal{S} sends \mathbf{V}''_{π} , \mathbf{V}'_{π} to \mathcal{R} .
3. For i -th row and j -column in \mathbf{V}''_{π} and \mathbf{V}'_{π} , if $v''_{\pi(ij)} = v'^a_{\pi(ij)}$, \mathcal{R} sets $b_{ij} = 1$, else, $b_{ij} = 0$. \mathcal{R} sets a bit matrix $\mathbf{B} = [b_{ij}]$, $i \in [\alpha]$, $j \in [m]$.

Fig. 11: Construction of pm-PEQT based on DDH

5.2 pm-PEQT based on DDH

The second construction of pm-PEQT is based on DDH as described in Figure 11. For notational convenience, we use $\text{pm-PEQT}_{\text{pub}}$ to denote DDH-based pm-PEQT, since the main computational costs are public-key operations.

Theorem 3. *The construction of Figure 11 securely implements functionality $\mathcal{F}_{\text{pm-PEQT}}$ based on DDH in the random oracle model, in the presence of semi-honest adversaries.*

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

Corrupt Sender. $\text{Sim}_{\mathcal{S}}(\mathbf{R}', \pi = (\pi_c, \pi_r))$ simulates the view of corrupt \mathcal{S} as follows: it chooses random group elements v_{ij} , $i \in [\alpha]$, $j \in [m]$ to simulate the view. We argue that the outputs of $\text{Sim}_{\mathcal{S}}$ are indistinguishable from the real view of \mathcal{S} by the following hybrids:

Hyb₀: \mathcal{S} 's view in the real protocol consists of $H(r_{ij})^a$, $i \in [\alpha]$, $j \in [m]$, where $a \leftarrow \mathbb{Z}_q$.

Hyb₁: Same as **Hyb₀** except that $\text{Sim}_{\mathcal{S}}$ chooses random group elements v_{ij} , $i \in [\alpha]$, $j \in [m]$ instead of $H(r_{ij})^a$, $i \in [\alpha]$, $j \in [m]$, where $a \leftarrow \mathbb{Z}_q$. The hybrid is the view output by $\text{Sim}_{\mathcal{S}}$.

We argue that the views in **Hyb₀** and **Hyb₁** are computationally indistinguishable. Let \mathcal{A} be a probabilistic polynomial-time (PPT) adversary against the DDH assumption. Given the DDH challenge $g^x, g^{y_{ij}}, g^{z_{ij}}$, where $x, y_{ij} \leftarrow \mathbb{Z}_q$, \mathcal{A} is asked to distinguish if $z_{ij} = x \cdot y_{ij}$ or random values. \mathcal{A} implicitly sets $a = x$, and simulates (with the knowledge of \mathbf{R}) the view as below:

- RO queries: $\text{Sim}_{\mathcal{S}}$ honestly emulates random oracle (RO) H . For queries r_{ij} , if $r_{ij} \notin \mathbf{R}$, it picks a random group element to assign $H(r_{ij})$, otherwise, it assigns $H(r_{ij}) = g^{y_{ij}}$.
- Outputs $g^{z_{ij}}$, $i \in [\alpha]$, $j \in [m]$.

Clearly, if $z_{ij} = x \cdot y_{ij}$, \mathcal{A} simulates Hyb_0 . Else, it simulates Hyb_1 (without the knowledge of \mathbf{R}), because it responds to all RO queries with random group elements without knowing that the inputs belong to \mathbf{R} or not. Therefore, the outputs of $\text{Sim}_{\mathcal{S}}$ are computationally indistinguishable from the real view based on the DDH assumption.

Corrupt Receiver. $\text{Sim}_{\mathcal{R}}(\mathbf{R}, \mathbf{B} = [b_{ij}])$ simulates the view of corrupt \mathcal{R} as follows: $\text{Sim}_{\mathcal{R}}$ chooses $a \leftarrow \mathbb{Z}_q$ randomly and simulates the first round message as real protocol. For $b_{ij} = 0$, $i \in [\alpha]$, $j \in [m]$, it chooses random group elements v_{ij} and u_{ij} to simulate the view. For $b_{ij} \neq 0$, $i \in [\alpha]$, $j \in [m]$, it chooses random group elements v_{ij} and sets $u_{ij} = v_{ij}^a$ to simulate the view.

We argue that the outputs of $\text{Sim}_{\mathcal{R}}$ are indistinguishable from the real view of \mathcal{R} by the following hybrids:

Hyb₀: \mathcal{R} 's view in the real protocol consists of $H(r'_{\pi(ij)})^b$ and $H(r_{\pi(ij)})^{ab}$, $i \in [\alpha]$, $j \in [m]$, where $a, b \leftarrow \mathbb{Z}_q$.

Hyb₁: Same as Hyb_0 except that for $b_{ij} = 0$, that is $r_{\pi(ij)} \neq r'_{\pi(ij)}$, $\text{Sim}_{\mathcal{R}}$ chooses random group elements v_{ij} and u_{ij} instead of $H(r'_{\pi(ij)})^b$ and $H(r_{\pi(ij)})^{ab}$.

Hyb₂: Same as Hyb_1 except that for $b_{ij} = 1$, that is $r_{\pi(ij)} = r'_{\pi(ij)}$, $\text{Sim}_{\mathcal{R}}$ chooses random group elements v_{ij} and sets $u_{ij} = v_{ij}^a$, $i \in [\alpha]$, $j \in [m]$ instead of $H(r'_{\pi(ij)})^b$ and $H(r_{\pi(ij)})^{ab}$. The hybrid is the view output by $\text{Sim}_{\mathcal{R}}$.

We argue that the view in Hyb_0 and Hyb_1 are computationally indistinguishable based on the DDH assumption. Given the DDH challenge $g^x, g^{y_{ij}}, g^{y'_{ij}}, g^{z_{ij}}, g^{z'_{ij}}$, where $x, y_{ij}, y'_{ij} \leftarrow \mathbb{Z}_q$, \mathcal{A} is asked to distinguish if $z_{ij} = x \cdot y_{ij}$, $z'_{ij} = x \cdot y'_{ij}$ or random values. \mathcal{A} implicitly sets $b = x$, and simulates (with the knowledge of \mathbf{R}' and π) the view as below:

- RO queries: $\text{Sim}_{\mathcal{R}}$ honestly emulates random oracle H . For queries r_{ij} and r'_{ij} , if $r_{ij} \notin \mathbf{R}$, $r'_{ij} \notin \mathbf{R}'$, it assigns $H(r_{\pi(ij)})$, $H(r'_{\pi(ij)})$ with random group elements. If $r_{ij} \in \mathbf{R}$, $r'_{ij} \in \mathbf{R}'$, it assigns $H(r_{ij}) = g^{a^{-1}y_{ij}}$, $H(r'_{ij}) = g^{y'_{ij}}$.
- Outputs $g^{z_{ij}}, g^{z'_{ij}}$.

Clearly, if $z_{ij} = x \cdot y_{ij}$, $z'_{ij} = x \cdot y'_{ij}$, \mathcal{A} simulates Hyb_0 . Else, it simulates Hyb_1 . In the Hyb_1 , $\text{Sim}_{\mathcal{R}}$ needs not to know the \mathbf{R}' and π in these positions with $b_{ij} = 0$, because in these positions, it responds to all random oracle queries with random group elements.

We argue that the view in Hyb_1 and Hyb_2 are computationally indistinguishable based on the DDH assumption. Given the DDH challenge $g^x, g^{y_{ij}}, g^{z_{ij}}$ where $x, y_{ij} \leftarrow \mathbb{Z}_q$, \mathcal{A} is asked to distinguish if $z_{ij} = x \cdot y_{ij}$ or random values. \mathcal{A} implicitly sets $b = x$, and simulates (with the knowledge of \mathbf{R}' and π for all positions with $b_{ij} = 1$) the view as below:

- RO queries: $\text{Sim}_{\mathcal{R}}$ honestly emulates random oracle H . For queries r_{ij} , r'_{ij} , if $r_{ij} \notin \mathbf{R}$, $r'_{ij} \notin \mathbf{R}'$, it assigns $H(r_{ij})$, $H(r'_{ij})$ with random group elements. If $r_{ij} = r'_{ij} \in \mathbf{R}$, it assigns $H(r_{ij}) = H(r'_{ij}) = g^{y_{ij}}$.
- Outputs $g^{a \cdot z_{ij}}, g^{z_{ij}}$.

Clearly, if $z_{ij} = x \cdot y_{ij}$, \mathcal{A} simulates Hyb_1 . Else, it simulates Hyb_2 . In the Hyb_2 , $\text{Sim}_{\mathcal{R}}$ needs not to know the \mathbf{R}' and π in these positions with $b_{ij} = 1$, because in these positions, it responds to all random oracle queries with random group elements. Therefore, the outputs of $\text{Sim}_{\mathcal{R}}$ are computationally indistinguishable from the real view based on the DDH assumption.

Protocols	Communication	Computation	Post-quantum security
pm-PEQT _{pub}	$O(m)$	$O(m)$	\times
pm-PEQT _{sym}	$O(m \log m)$	$O(m \log m)$	\checkmark

Table 2: Comparisons of two pm-PEQT constructions. m denotes the number of elements in the matrix.

Comparisons of two pm-PEQT constructions. pm-PEQT_{pub} based on DDH assumption is not post-quantum secure, while pm-PEQT_{sym} based on symmetric cryptographic primitives is potentially post-quantum secure. The comparison of the two pm-PEQT constructions is given in Table 2.

6 Full PSU Protocol

In this section, we detail our full PSU protocol in Figure 12, given a secure fully homomorphic encryption scheme with circuit privacy. The ideal functionality of this protocol is given in Figure 4.

Correctness. The correctness of our PSU protocol is conditioned on the hashing succeeding, which happens with overwhelming probability $1 - 2^{-\lambda}$.

Theorem 4. *The protocol in Figure 12, is a secure protocol for \mathcal{F}_{PSU} in the $(\mathcal{F}_{pm-PEQT}, \mathcal{F}_{OT})$ -hybrid model, in the presence of semi-honest adversaries, provided that the fully homomorphic encryption scheme is IND-CPA secure with circuit privacy.*

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

Corrupt Sender. $\text{Sim}_{\mathcal{S}}(X)$ simulates the view of corrupt \mathcal{S} as follows. $\text{Sim}_{\mathcal{S}}$ hashes X into X_c as the real protocol, and encrypts random values in place of the ciphertexts in step 5. Then it decrypts the ciphertexts as \mathbf{R}' and chooses randomly permutation $\pi = (\pi_c, \pi_r)$. It invokes $\text{Sim}_{pm-PEQT}^{\mathcal{S}}(\mathbf{R}', \pi)$ and $\text{Sim}_{OT}^{\mathcal{S}}(\perp, X_c[\pi_c(j)])$, $j \in [m_c]$ appends the output to the view. Now we argue that the view output by $\text{Sim}_{\mathcal{S}}$ is indistinguishable from the real one. The plaintexts are randomized in the real view which is indistinguishable from the random values in the simulated view. The FHE satisfies the circuit privacy which hides the computational circuit. The views of the underlying pm-PEQT and OT simulator are indistinguishable. Thus, the simulation is indistinguishable from the real view.

Corrupt Receiver. $\text{Sim}_{\mathcal{R}}(Y, X \cup Y)$ simulates the view of corrupt receiver as follows: $\text{Sim}_{\mathcal{R}}$ encrypts random value in place of the ciphertexts in step 4. It chooses the random matrix \mathbf{R} in step 3. $\text{Sim}_{\mathcal{R}}$ computes $\hat{X} = (X \cup Y) \setminus Y$ and pads \hat{X} with \perp to m_c items and permutes these items randomly. For all items in \hat{X} , if $\hat{x}_i \neq \perp$, it sets $b_i = 1$, else $b_i = 0$. And then it generates $\mathbf{B}_{\alpha \times m_c}$, for all columns \mathbf{b}_i , if $b_i = 1$, all items in \mathbf{b}_i are set to 0, else, one random position in \mathbf{b}_i is set to 1 and all other positions are set to 0. $\text{Sim}_{\mathcal{R}}$ invokes $\text{Sim}_{pm-PEQT}^{\mathcal{R}}(\mathbf{R}, \mathbf{B})$ appends the output to the view. Then, for all $i \in [m_c]$, it invokes $\text{Sim}_{OT}^{\mathcal{R}}(b_i, \hat{x}_i)$ and appends the output to the view.

The view generated by $\text{Sim}_{\mathcal{R}}$ is indistinguishable from a real view of \mathcal{R} by the following hybrids:

Hyb₀: \mathcal{R} 's view in the real protocol.

Hyb₁: Same as **Hyb₀** except that the ciphertexts are replaced by encrypting random values generated by $\text{Sim}_{\mathcal{R}}$. Since the fully homomorphic encryption scheme is IND-CPA secure, the simulation is indistinguishable from the real view.

Hyb₂: Same as **Hyb₁** except that the output of $\mathcal{F}_{pm-PEQT}$ is replaced by \mathbf{B} generated by $\text{Sim}_{\mathcal{R}}$, and $\text{Sim}_{\mathcal{R}}$ runs the $\mathcal{F}_{pm-PEQT}$ simulator to produce the simulated view for \mathcal{R} . The security of the pm-PEQT protocol guarantees the view in simulation is computationally indistinguishable from the view in the real protocol.

Hyb₃: Same as **Hyb₂** except that $\text{Sim}_{\mathcal{R}}$ runs the \mathcal{F}_{OT} simulator to produce the simulated view for \mathcal{R} . The security of OT protocol guarantees the view in simulation is computationally indistinguishable from the view in the real protocol. The hybrid is the view output by $\text{Sim}_{\mathcal{R}}$.

Note that, the steps 1-6 in our full PSU protocol can be seen as the opposite case of unbalanced PSI [CLR17]. Therefore, same optimizations [CLR17], such as batching, windowing, partitioning, and modulus switching, are suitable for the steps 1-6 to significantly reduce the depth of the homomorphic circuit and improve the efficiency. We review the optimizations in Appendix A.

Communication. The analysis of the communication complexity is provided as follows. In steps 1-6, the communication requires $O(m \log n)$; the communication of pm-PEQT requires $O(m \log m)$ (based on Permute

Input: The receiver inputs set $Y \subset \{0,1\}^*$ of size $n = |Y|$ and the sender inputs set $X \subset \{0,1\}^*$ of size $m = |X|$, where $m \ll n$, m and n are public.

Output: The receiver outputs $X \cup Y$. The sender outputs \perp .

- (1) **[Setup]** \mathcal{R} and \mathcal{S} agree on the hashing, FHE, mp-PEQT and OT parameters.
- (2) **[Hashing]** \mathcal{S} hashes X into table X_c by Cuckoo hashing, where X_c consists of m_c bins and each bin has only one item. \mathcal{R} uses the same hash functions to hash Y into table $\mathbf{Y}_{B \times m_c}$ by simple hashing, where $\mathbf{Y}_{B \times m_c}$ consists of m_c bins and each bin has B items.
- (3) **[Pre-process Y]**
 - (a) **[Partitioning]** \mathcal{R} partitions $\mathbf{Y}_{B \times m_c}$ by rows into α subtables $\mathbf{Y}_1, \dots, \mathbf{Y}_\alpha$. Each subtable has $B' = B/\alpha$ rows and m_c columns.
 - (b) **[Computing coefficients]** For j -th columns of i -th subtable $\mathbf{y}_{i,j} = [y_{i,j,1}, \dots, y_{i,j,B'}]^T$, $i \in [\alpha], j \in [m_c]$, \mathcal{R} computes the coefficients of the polynomial $f_{i,j}(y) = \prod_{k=1}^{B'} (y - y_{i,j,k}) = a'_{i,j,0} + a_{i,j,1}y + \dots + a_{i,j,B'}y^{B'}$. \mathcal{R} computes the coefficient matrix \mathbf{A} as follows: \mathcal{R} chooses a random matrix $\mathbf{R}_{\alpha \times m_c} = [r_{i,j}]$, and sets j -th columns of i -th subtable $\mathbf{A}_{i,j} = [a_{i,j,0}, a_{i,j,1}, \dots, a_{i,j,B'}]^T$, $i \in [\alpha], j \in [m_c]$, where $a_{i,j,0} = a'_{i,j,0} + r_{i,j}$.
 - (c) **[Batching]** For each subtable obtained from the previous step, \mathcal{R} interprets each of its row as a vector of length m_c with elements in \mathbb{Z}_t . Then \mathcal{R} batches each vector into $\beta = m_c/\gamma$ plaintext polynomials. As a result, each row of i -th subtable \mathbf{A}_i is transformed into β polynomials denoted $\hat{\mathbf{A}}_{i,j}$, $i \in [\alpha]$, $j \in [\beta]$.
- (4) **[Encrypt X]**
 - (a) **[Batching]** \mathcal{S} interprets X_c as a vector of length m_c with items in \mathbb{Z}_t , and batches this vector into $\beta = m_c/\gamma$ plaintext polynomials $\hat{X}_1, \dots, \hat{X}_\beta$.
 - (b) **[Windowing]** For each batched plaintext polynomial \hat{X} , \mathcal{S} computes the component-wise $i \cdot 2^j$ -th powers $\hat{X}^{i \cdot 2^j}$, for $i \in [2^l - 1]$ and $0 \leq j \leq \lceil \log_2(B')/l \rceil$.
 - (c) **[Encrypt]** \mathcal{S} uses FHE scheme to encrypt each such power, obtaining β collections of ciphertexts $\mathbf{C}_j, j \in [\beta]$, and each collection consists of the ciphertexts $[c_{i,j}], i \in [2^l - 1], 0 \leq j \leq \lceil \log_2(B')/l \rceil$. \mathcal{S} sends these ciphertexts to \mathcal{R} .
- (5) **[Computation]**
 - (a) **[Homomorphically compute encryptions of all powers]** For each collection $\mathbf{C}_j, j \in [\beta]$, \mathcal{R} homomorphically computes encryptions of all powers $\mathbf{C}_j = [\mathbf{c}_{j,0}, \dots, \mathbf{c}_{j,B'}]$, where $\mathbf{c}_{j,k}, 0 \leq k \leq B'$ is a ciphertext encrypted \hat{X}_j^k .
 - (b) **[Homomorphically evaluate the dot product]** \mathcal{R} homomorphically evaluates $\mathbf{C}'_{i,j} = \mathbf{C}_j \hat{\mathbf{A}}_{i,j}, i \in [\alpha], j \in [\beta]$, performs modulus switching on $\mathbf{C}'_{i,j}$ to reduce sizes, and sends the ciphertexts to \mathcal{S} .
- (6) **[Decrypt]** \mathcal{S} decrypts the ciphertexts and concatenates the results into the matrix $\mathbf{R}'_{\alpha \times m_c}$.
- (7) **[pm-PEQT]** \mathcal{R} inputs the matrix $\mathbf{R}_{\alpha \times m_c}$, and \mathcal{S} inputs the permutation $\pi = (\pi_c, \pi_r)$ and the matrix $\mathbf{R}'_{\alpha \times m_c}$. Both parties invoke the pm-PEQT functionality. As a result, \mathcal{R} gets a bit matrix $\mathbf{B}_{\alpha \times m_c} = [b_{ij}]$, where $b_{ij} = 1$ if $r_{\pi(ij)} = r'_{\pi(ij)}$, and $b_{ij} = 0$ otherwise.
- (8) **[Output]** \mathcal{R} sets a bit vector $\mathbf{b} = [b_j], j \in [m_c]$, where $b_j = 1$ if for all $i \in [\alpha]$, $b_{ij} = 0$ and $b_j = 0$ otherwise. Then, \mathcal{R} and \mathcal{S} invoke the OT functionality, in which \mathcal{R} inputs $b_j, j \in [m_c]$ and \mathcal{S} inputs $(\perp, X_c[\pi_c(j)])$. If $b_j = 1$, \mathcal{R} gets $X_c[\pi_c(j)]$, else, gets \perp . Finally, \mathcal{R} outputs the set union $Y \cup \{X_c[\pi_c(j)]\}$, for all $j \in [m_c]$.

Fig. 12: Full PSU protocol

+ Share and mp-OPRF) or $O(m)$ (based on DDH) since we omit the parameter α which is used to make a trade-off between the computation and communication of our PSU; the communication of OT requires $O(m)$. In summary, the communication of our PSU is $O(m \log n)$.

We describe some further optimization ideas of our protocol as follows.

Offline/online. Following [CLR17, CHLR18, CMdG+21], the pre-processing of the receiver in our PSU can be done entirely offline without involving the sender. Specifically, given an upper bound on the sender’s set size, the receiver can locally choose parameters and perform the pre-processing. Upon learning the sender’s actual set size, the receiver sends the parameters to the sender, and the sender pads the same dummy items like \perp which are known to both parties.

OPRF pre-processing. The security of our full protocol is proven based on FHE with circuit privacy following [CLR17]. This leads to performing a noise flooding operation on the result ciphertexts, as is necessary in [CLR17]. Following [CHLR18, CMdG+21], we can use an OPRF to compute the items on both sides before engaging in the PSU, which prevents the sender from learning anything about the original items. Thus, our PSU can be proven secure without circuit privacy and utilize more efficient FHE parameters as [CHLR18, CMdG+21], which improves our performance and adds flexibility to the parametrization.

7 Implementation and Performance

In this section, we experimentally evaluate our two PSU protocols:

- PSU_{sym} : PSU protocol based on FHE, OT, and pm-PEQT, where pm-PEQT is built from Permute + Share and mp-OPRF.
- PSU_{pub} : PSU protocol based on FHE, OT, and DDH-based pm-PEQT.

We first give the experimental environment, and then compare our protocols with the state-of-the-art works in terms of communication and runtime in different network environments. Our source code is available upon request.

7.1 Experimental Setup

We run our experiments on a single Intel Core i7-11700 CPU @ 2.50GHz with 16 threads and 16GB of RAM. We simulate network latency and bandwidth by using the Linux `tc` command. Specifically, we consider the following LAN setting, where the two parties are connected via a local host with 10Gbps throughput, and a 0.2ms round-trip time (RTT). We also consider three WAN settings with 100Mbps, 10Mbps and 1Mbps bandwidth, each with an 80ms RTT.

7.2 Implementation Details

We use the FHE scheme in [FV12], Permute + Share in [MS13, CGP20], mp-OPRF in [CM20] and OT extension in [ALSZ13]. Following [JSZ+22, ZCL+23], we set the computational security parameter $\kappa = 128$ and the statistical security parameter $\lambda = 40$. Our implementation is written in C++. The following libraries are used in our implementation.

- FHE: SEAL <https://github.com/microsoft/SEAL>
and APSI <https://github.com/microsoft/APSI>
- Permute + Share: <https://github.com/dujjiajun/PSU>
- mp-OPRF: <https://github.com/peihanmiao/OPRF-PSI>
- OT: <https://github.com/osu-crypto/libOTe>

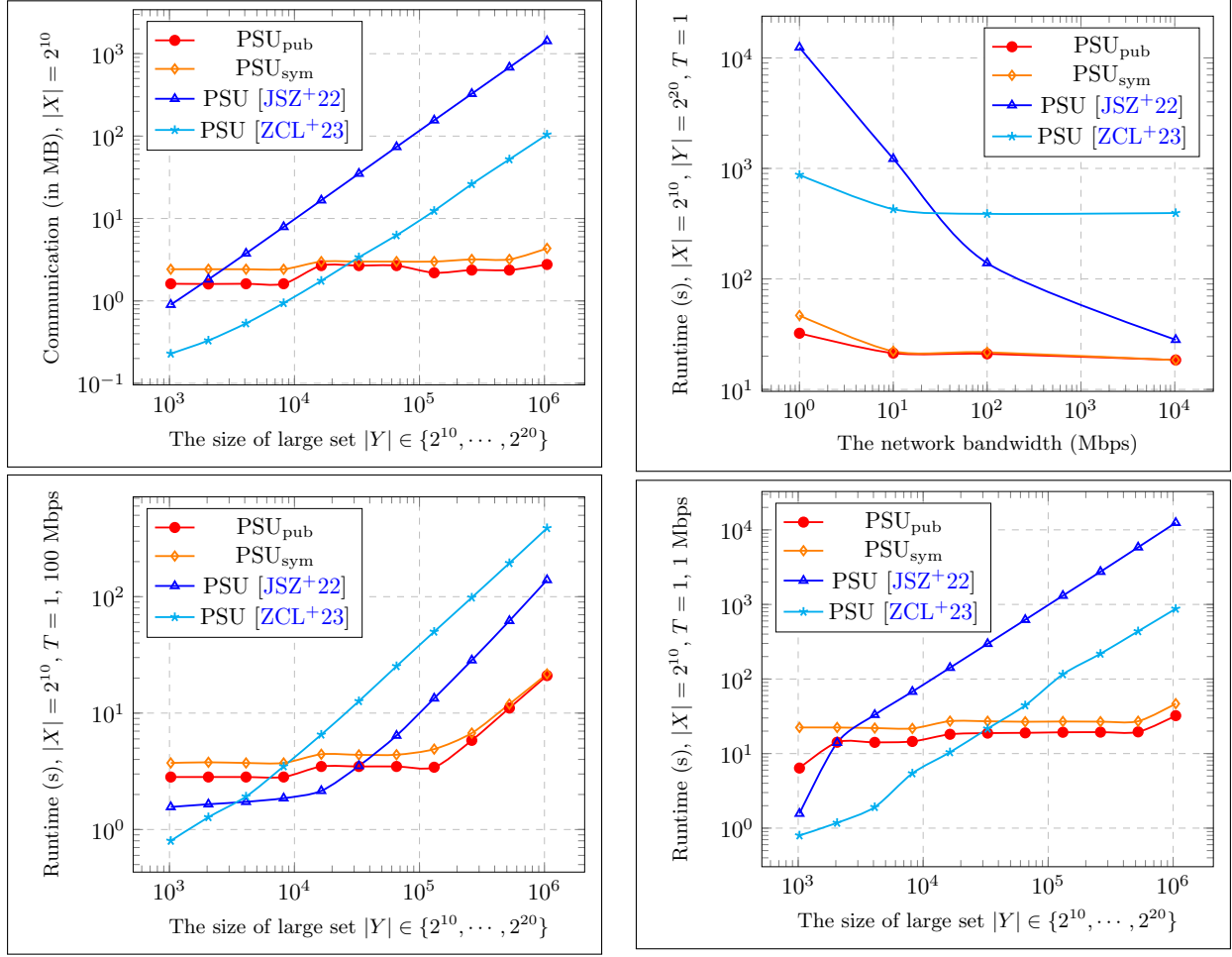


Fig.13: Comparisons of communication (in MB) and runtime (in seconds) between PSU [JSZ+22], PSU [ZCL+23], PSU_{pub} and PSU_{sym}. Both x and y -axis are in log scale. The top left figure shows the communication cost increases as the large set size increases. The top right figure shows the runtime decreases as the bandwidth increases. The bottom two figures show the runtime increases as the large set size increases.

7.3 Performance Comparisons

In this section, we compare PSU_{sym} and PSU_{pub} with PSU [JSZ+22, ZCL+23] in terms of runtime and communication, and the results are reported in Table 3 and Figure 13.

We stress that all reported costs are collected in the same environment. For comparisons of other works [JSZ+22, ZCL+23], we use the parameters recommended in their open-source code and fix the item length to 128-bit.

- PSU [JSZ+22]: <https://github.com/dujiajun/PSU>
- PSU [ZCL+23]: <https://github.com/alibaba-edu/mpc4j>

Communication comparison. Our PSU protocol can achieve the lowest communication among all protocols [JSZ+22, ZCL+23] in the unbalanced setting. We provide an analysis of the asymptotic border of communication beyond which our protocols perform better than [ZCL+23]. The detailed analysis is as below: the communication complexity of our PSU is $O(m \log n)$, while the communication complexity of the

state-of-the-art PSU protocol due to Zhang et al. [ZCL+23] is $O(n + m)$. Therefore, in the asymptotic sense the performance border between our protocol and Zhang et al.’s protocol can be deduced from the inequality $m < \frac{n}{\log n - 1}$. Our experimental results shown in Table 3 is consistent with our asymptotic analysis: when ($n \geq 2^{15}$, $m = 2^{10}$), our PSU protocols are superior than [ZCL+23]. Especially, as shown in Figure 13, the larger difference between two set sizes, the better our protocols perform. For set sizes ($|X| = 2^{10}$, $|Y| = 2^{20}$), the communication of our PSU_{pub} requires 2.767 MB, which is about $37\times$ lower than PSU [ZCL+23] requiring 104.28 MB, about $515\times$ lower than PSU [JSZ+22] requiring 1426.855MB.

Runtime comparison. Our PSU_{sym} and PSU_{pub} are faster than PSU [JSZ+22, ZCL+23] in the unbalanced case depending on network environments. As shown in Figure 13, the larger difference between the two set sizes, the better our protocols perform. For set sizes ($|X| = 2^{10}$, $|Y| = 2^{20}$) with $T = 1$ thread in LAN setting, the runtime of our PSU_{pub} requires 18.465 seconds, while PSU [ZCL+23] requires 394.92 seconds, about $21\times$ improvement, PSU [JSZ+22] requires 28.137 seconds, about $1.5\times$ improvement. The performance of our protocols improves significantly in the case of low bandwidth. For set sizes ($|X| = 2^{10}$, $|Y| = 2^{20}$) with $T = 8$ thread in 1 Mbps bandwidth, our PSU_{pub} requires 24.751 seconds, while PSU [ZCL+23] requires 874.854 seconds, about $35\times$ improvement, PSU [JSZ+22] requires 12115.414 seconds, about $489\times$ improvement.

7.4 Scalability and Parallelizability

We demonstrate the scalability and parallelizability of our protocol by evaluating it on large set sizes $n = \{2^{10}, 2^{11}, 2^{12}, 2^{13}, 2^{14}, 2^{15}, 2^{16}, 2^{17}, 2^{18}, 2^{19}, 2^{20}\}$ and fixed small set size $m = 2^{10}$. We run each party in parallel with $T \in \{1, 2, 8\}$ threads and give the performance of our PSU_{pub} in Table 4, showing running time in both LAN/WAN settings.

Our protocol indeed scales well in the unbalanced setting. In single-thread, for set size $(2^{10}, 2^{10})$, our protocol requires 0.916 seconds; for set size $(2^{10}, 2^{15})$, requires 1.272 seconds; for set size $(2^{10}, 2^{20})$, requires 18.465 seconds. Increasing the number of threads to $T = 8$ runs the $(2^{10}, 2^{20})$ instance in 6.71 seconds. Benchmarking our implementation in the WAN setting, our protocol also scales well due to the low communication cost (for $(2^{10}, 2^{20})$, our protocol needs 2.767 MB of communication which is about $37\times$ lower than PSU [ZCL+23] and about $515\times$ lower than PSU [JSZ+22]).

Our protocol is very amenable to parallelization. Specifically, following [CLR17] divides (n, m) -PSI into many $(B', 1)$ -PSI, (m, n) -PSU is divided into many $(1, B')$ -PSU. Thus, on the side of the large set, the receiver can compute all polynomials and ciphertexts in parallel. For example, when increasing T from 1 to 8, our protocol shows a factor of $2.75\times$ improvement as the running time reduces from 18.465 seconds to 6.71 seconds for set size $(2^{10}, 2^{20})$.

We presents the ratio between the runtime of the single thread and 8 threads in the last row of Table 4. Our protocol yields a better speedup when the large set size become larger. For large set size of $n = 2^{10}$, the protocol achieves a moderate speedup of about 1.62. When considering the larger size $n = 2^{20}$, the speedup of about 2.75 is obtained at 8 threads.

References

- [ALSZ13] Gilad Asharov, Yehuda Lindell, Thomas Schneider, and Michael Zohner. More efficient oblivious transfer and extensions for faster secure computation. In Ahmad-Reza Sadeghi, Virgil D. Gligor, and Moti Yung, editors, *2013 ACM SIGSAC Conference on Computer and Communications Security, CCS'13, Berlin, Germany, November 4-8, 2013*, pages 535–548. ACM, 2013.
- [BA12] Marina Blanton and Everaldo Aguiar. Private and oblivious set and multiset operations. In *ASIACCS*, 2012.
- [BGH13] Zvika Brakerski, Craig Gentry, and Shai Halevi. Packed ciphertexts in lwe-based homomorphic encryption. In *Public-Key Cryptography - PKC 2013*, pages 1–13, 2013.
- [BGV12] Zvika Brakerski, Craig Gentry, and Vinod Vaikuntanathan. (leveled) fully homomorphic encryption without bootstrapping. In *Innovations in Theoretical Computer Science 2012*, pages 309–325, 2012.

- [BPMW16] Florian Bourse, Rafaël Del Pino, Michele Minelli, and Hoeteck Wee. FHE circuit privacy almost for free. In *CRYPTO*, 2016.
- [BSMD10] Martin Burkhart, Mario Strasser, Dilip Many, and Xenofontas A. Dimitropoulos. SEPIA: privacy-preserving aggregation of multi-domain network events and statistics. In *19th USENIX Security Symposium*, pages 223–240, 2010.
- [CCS18] Andrea Cerulli, Emiliano De Cristofaro, and Claudio Soriente. Nothing refreshes like a repsi: Reactive private set intersection. In *Applied Cryptography and Network Security ACNS 2018*, pages 280–300, 2018.
- [CGP20] Melissa Chase, Esha Ghosh, and Oxana Poburinnaya. Secret-shared shuffle. In *Advances in Cryptology - ASIACRYPT 2020*, pages 342–372, 2020.
- [CHLR18] Hao Chen, Zhicong Huang, Kim Laine, and Peter Rindal. Labeled PSI from fully homomorphic encryption with malicious security. In *ACM SIGSAC Conference on Computer and Communications Security*, pages 1223–1237, 2018.
- [CLP17] Hao Chen, Kim Laine, and Rachel Player. Simple encrypted arithmetic library - SEAL v2.1. In *Financial Cryptography and Data Security - FC 2017*, pages 3–18, 2017.
- [CLR17] Hao Chen, Kim Laine, and Peter Rindal. Fast private set intersection from homomorphic encryption. In *ACM SIGSAC Conference on Computer and Communications Security*, pages 1243–1255, 2017.
- [CM20] Melissa Chase and Peihan Miao. Private set intersection in the internet setting from lightweight oblivious PRF. In *Advances in Cryptology - CRYPTO 2020*, pages 34–63, 2020.
- [CMdG⁺21] Kelong Cong, Radames Cruz Moreno, Mariana Botelho da Gama, Wei Dai, Iliia Iliashenko, Kim Laine, and Michael Rosenberg. Labeled PSI from homomorphic encryption with reduced computation and communication. In *ACM SIGSAC Conference on Computer and Communications Security*, pages 1135–1150, 2021.
- [CO18] Michele Ciampi and Claudio Orlandi. Combining private set-intersection with secure two-party computation. In *Security and Cryptography for Networks - 11th International Conference, SCN 2018*, pages 464–482, 2018.
- [CZZD22] Yu Chen, Min Zhang, Cong Zhang, and Minglang Dong. Private set operations from multi-query reverse private membership test. *Cryptology ePrint Archive*, Paper 2022/652, 2022.
- [DC17] Alex Davidson and Carlos Cid. An efficient toolkit for computing private set operations. In *Information Security and Privacy - 22nd Australasian Conference, ACISP 2017*, pages 261–278, 2017.
- [DCW13] Changyu Dong, Liqun Chen, and Zikai Wen. When private set intersection meets big data: an efficient and scalable protocol. In *2013 ACM SIGSAC Conference on Computer and Communications Security, CCS'13*, pages 789–800, 2013.
- [DM03] Luc Devroye and Pat Morin. Cuckoo hashing: Further analysis. *Inf. Process. Lett.*, 86(4):215–219, 2003.
- [FNO19] Brett Hemenway Falk, Daniel Noble, and Rafail Ostrovsky. Private set intersection with linear communication from general assumptions. In *Proceedings of the 18th ACM Workshop on Privacy in the Electronic Society*, pages 14–25, 2019.
- [FPSS03] Dimitris Fotakis, Rasmus Pagh, Peter Sanders, and Paul G. Spirakis. Space efficient hash tables with worst case constant access time. In *STACS 2003*, pages 271–282, 2003.
- [Fri07] Keith B. Frikken. Privacy-preserving set union. In *Applied Cryptography and Network Security, ACNS 2007*, pages 237–252, 2007.
- [FV12] Junfeng Fan and Frederik Vercauteren. Somewhat practical fully homomorphic encryption. *IACR Cryptol. ePrint Arch.*, page 144, 2012.
- [GDL⁺16] Ran Gilad-Bachrach, Nathan Dowlin, Kim Laine, Kristin E. Lauter, Michael Naehrig, and John Wernsing. Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy. In *Proceedings of the 33rd International Conference on Machine Learning, ICML*, pages 201–210, 2016.
- [GHS12] Craig Gentry, Shai Halevi, and Nigel P. Smart. Homomorphic evaluation of the AES circuit. In *Advances in Cryptology - CRYPTO 2012*, pages 850–867, 2012.
- [GMR⁺21] Gayathri Garimella, Payman Mohassel, Mike Rosulek, Saeed Sadeghian, and Jaspal Singh. Private set operations from oblivious switching. In *Public-Key Cryptography - PKC 2021*, pages 591–617, 2021.
- [HEK12] Yan Huang, David Evans, and Jonathan Katz. Private set intersection: Are garbled circuits better than custom protocols? In *19th Annual Network and Distributed System Security Symposium, 2012*, 2012.

- [HLS⁺16] Kyle Hogan, Noah Luther, Nabil Schear, Emily Shen, David Stott, Sophia Yakoubov, and Arkady Yerukhovich. Secure multiparty computation for cooperative cyber risk assessment. In *IEEE Cybersecurity Development, SecDev 2016*, pages 75–76, 2016.
- [HV17] Carmit Hazay and Muthuramakrishnan Venkatasubramanian. Scalable multi-party private set-intersection. In *Public-Key Cryptography - PKC 2017*, pages 175–203, 2017.
- [IKNP03] Yuval Ishai, Joe Kilian, Kobbi Nissim, and Erez Petrank. Extending oblivious transfers efficiently. In *Advances in Cryptology - CRYPTO 2003*, Lecture Notes in Computer Science, 2003.
- [JSZ⁺22] Yanxue Jia, Shi-Feng Sun, Hong-Sheng Zhou, Jiajun Du, and Dawu Gu. Shuffle-based private set union: Faster and more secure. *IACR Cryptol. ePrint Arch.*, page 157, 2022.
- [KLS⁺17] Ágnes Kiss, Jian Liu, Thomas Schneider, N. Asokan, and Benny Pinkas. Private set intersection for unequal set sizes with mobile applications. *Proc. Priv. Enhancing Technol.*, pages 177–197, 2017.
- [KMP⁺17] Vladimir Kolesnikov, Naor Matania, Benny Pinkas, Mike Rosulek, and Ni Trieu. Practical multi-party private set intersection from symmetric-key techniques. In *ACM SIGSAC Conference on Computer and Communications Security*, pages 1257–1272, 2017.
- [KRTW19] Vladimir Kolesnikov, Mike Rosulek, Ni Trieu, and Xiao Wang. Scalable private set union from symmetric-key techniques. In *Advances in Cryptology - ASIACRYPT 2019*, pages 636–666, 2019.
- [KS05] Lea Kissner and Dawn Xiaodong Song. Privacy-preserving set operations. In *Advances in Cryptology - CRYPTO 2005*, pages 241–257, 2005.
- [LV04] Arjen K. Lenstra and Tim Voss. Information security risk assessment, aggregation, and mitigation. In *Information Security and Privacy: 9th Australasian Conference, ACISP*, pages 391–401, 2004.
- [MR95] Rajeev Motwani and Prabhakar Raghavan. Randomized algorithms. Cambridge university press, 1995.
- [MS13] Payman Mohassel and Seyed Saeed Sadeghian. How to hide circuits in MPC an efficient framework for private function evaluation. In *Advances in Cryptology - EUROCRYPT 2013*, pages 557–574, 2013.
- [PR01] Rasmus Pagh and Flemming Friche Rodler. Cuckoo hashing. In *Algorithms - ESA 2001*, pages 121–133, 2001.
- [PRTY19] Benny Pinkas, Mike Rosulek, Ni Trieu, and Avishay Yanai. Spot-light: Lightweight private set intersection from sparse OT extension. In *Advances in Cryptology - CRYPTO 2019*, pages 401–431, 2019.
- [PSTY19] Benny Pinkas, Thomas Schneider, Oleksandr Tkachenko, and Avishay Yanai. Efficient circuit-based PSI with linear communication. In *Advances in Cryptology - EUROCRYPT 2019*, pages 122–153, 2019.
- [PSZ14] Benny Pinkas, Thomas Schneider, and Michael Zohner. Faster private set intersection based on OT extension. In *Proceedings of the 23rd USENIX Security Symposium*, pages 797–812, 2014.
- [PSZ18] Benny Pinkas, Thomas Schneider, and Michael Zohner. Scalable private set intersection based on OT extension. *ACM Trans. Priv. Secur.*, pages 7:1–7:35, 2018.
- [RA18] Amanda Cristina Davi Resende and Diego F. Aranha. Faster unbalanced private set intersection. pages 203–221, 2018.
- [Rab05] Michael O. Rabin. How to exchange secrets with oblivious transfer. *IACR Cryptol. ePrint Arch.*, page 187, 2005.
- [RMY20] Sivaramakrishnan Ramanathan, Jelena Mirkovic, and Minlan Yu. BLAG: improving the accuracy of blacklists. In *NDSS*. The Internet Society, 2020.
- [RR17] Peter Rindal and Mike Rosulek. Improved private set intersection against malicious adversaries. In *Advances in Cryptology - EUROCRYPT 2017*, pages 235–259, 2017.
- [SV14] Nigel P. Smart and Frederik Vercauteren. Fully homomorphic SIMD operations. *Des. Codes Cryptogr.*, pages 57–81, 2014.
- [ZCL⁺23] Cong Zhang, Yu Chen, Weiran Liu, Min Zhang, and Dongdai Lin. Optimal private set union from multi-query reverse private membership test. In *USENIX 2023*, 2023. <https://eprint.iacr.org/2022/358>.

A Optimization techniques

Our PSU considers the opposite case of [CLR17], where \mathcal{R} holds the large set and \mathcal{S} holds the small set. Thus, we can take advantage of the same optimization techniques of FHE following [CLR17], such as batching,

windowing, partitioning, and modulus switching, to significantly reduce the depth of the homomorphic circuit. We review the optimization techniques as follows.

Batching. Batching is a well-known and powerful technique in fully homomorphic encryption to enable Single Instruction, Multiple Data (SIMD) operations on ciphertexts [GHS12, BGH13, SV14, CLP17, GDL⁺16]. The batching technique allows the receiver to operate on γ items from the sender simultaneously, resulting in γ -fold improvement in both the computation and communication. As an example, The sender groups its items into vectors of length γ , encrypts them, and sends m/γ ciphertexts to the receiver. Upon receiving each ciphertext c_i , the receiver samples a vector $\mathbf{r}_i = (r_{i1}, \dots, r_{i\gamma}) \in (\mathbb{Z}_t \setminus \{0\})^n$ at random, homomorphically computes $r_i + \prod_{y \in Y} (c_i - y)$, and sends it back to the sender. Note that these modifications do not affect correctness or security, since the exact same proof can be applied to per each vector coefficient.

Windowing. We use the windowing technique [CLR17] to lower the depth of the circuit. In our PSU, the receiver needs to evaluate the sender’s encrypted data. If the receiver only has an encryption $c \leftarrow \text{FHE.Enc}(x)$, it samples a random r in $\mathbb{Z}_t \setminus \{0\}$ and homomorphically computes $r + \prod_{y_i \in Y} (c - y_i)$. The receiver computes at worst the product x^n , which requires a circuit of depth $\lceil \log_2 n \rceil$. To see this, we write $r + \prod_{y_i \in Y} (x - y_i) = r + a_0 + a_1x + \dots + a_{n-1}x^{n-1} + x^n$. If the sender sends encryptions of extra powers of x , the receiver could use these powers to evaluate the same computation with a much lower-depth circuit. More precisely, for a window size of l bits, the sender computes and sends $c_{ij} = \text{FHE.Enc}(x^{i \cdot 2^{lj}})$ to the receiver for all $1 \leq i \leq 2^l - 1$, $0 \leq j \leq \lfloor \log_2(n)/l \rfloor$. For example, when $l = 1$, the receiver sends encryptions of $x, x^2, \dots, x^{2^{\lfloor \log_2 n \rfloor}}$. This technique results in a significant reduction in the circuit depth.

Partitioning. Another way to reduce circuit depth is to let the receiver partition its set into α subsets [CLR17, CHLR18]. In our PSU, the receiver needs to compute encryptions of all powers x, \dots, x^n for each of the sender’s items x . With partitioning, the receiver only needs to compute encryptions of $x, \dots, x^{n/\alpha}$, which can be reused for each of the α partitions.

Modulus switching. We employ modulus switching to effectively reduce the size of the response ciphertexts as [CLR17, CHLR18, BGV12]. Modulus switching is a well-known operation in lattice-based fully homomorphic encryption schemes. It is a public operation, which transforms a ciphertext with encryption parameter q into a ciphertext encrypting the same plaintext, but with a smaller parameter $q' < q$. Note that the security of the protocol is trivially preserved as long as the smaller modulus q' is determined at setup.