

# Communication-Efficient and Privacy-Preserving Protocol for Computing Over-Threshold Set-Union

Xuhui Gong, Qiang-sheng Hua<sup>(✉)</sup>, and Hai Jin

National Engineering Research Center–Big Data Technology and System Lab, Key  
Laboratory of Services Computing Technology and System, Key Laboratory of  
Cluster and Grid Computing, School of Computer Science and Technology, Huazhong  
University of Science and Technology, Wuhan 430074, P.R. China  
qshua@hust.edu.cn

**Abstract.** In a variety of applications, the data items of multiple participants are collected and analyzed, and meanwhile the participants' privacy needs to be protected. This paper studies an over-threshold data aggregation problem, i.e., over-threshold set-union. In our model, we assume there are  $n$  participants, an untrusted data aggregator and a proxy, and each participant has a sensitive data item. The over-threshold set-union is normally defined as follows: given a threshold  $t$ , the aggregator only learns the data items which appear at least  $t$  times in the union of data items of all participants without learning anything else. Existing solutions either suffer from high communication cost or leak the multiplicity information of data items. In order to handle this defect, we present an efficient protocol in the honest-but-curious model by leveraging threshold secret sharing and dual pairing vector spaces. We prove that the proposed protocol not only has  $O(n(\log n)^2)$  communication complexity which nearly matches the lower bound  $\Omega(n/\log n)$  but also protects the data privacy.

**Keywords:** Over-threshold set-union · Data aggregation · Secure multi-party computation · Privacy

## 1 Introduction

The privacy-preserving data aggregation problem is of particular interest in many important applications. Most applications need to collect participants' data and compute some algebraic statistics. However, the data items of participants usually include sensitive information and the participants are reluctant to leak the data privacy to anyone. For example [10], in privacy-preserving distributed network monitoring, every node monitors local behavior and a service provider is responsible for a global anomalous monitoring task, and when the number of local anomalous behaviors exceeds a predetermined system threshold, the service provider identifies the anomaly.

There are two major challenges for the privacy-preserving data aggregation problem. One is communication efficiency. In many practical systems, monitoring nodes usually have tight resource constraints. Thus, minimizing communication overhead is a pursued goal of related researchers for designing the protocol. Another issue is privacy-preserving. After a protocol has been executed, the result of the computation is leaked and no adversary can learn more extra information than what can be inferred from the result.

Many types of research have been proposed on the privacy-preserving data aggregation problem. For example, linear aggregation function-SUM [4], [8], [11], and non-linear aggregation functions-MIN/MAX [11] and top- $k$  [5]. In this paper, we study a particular privacy-preserving data aggregation problem, i.e., the over-threshold set-union. This problem can be used in many application scenarios, such as the privacy-preserving distributed network monitoring [10]. As far as we know, there are some works [10], [5], [2], [9] that discuss the over-threshold set-union. In [10], the authors present a probabilistic protocol based on polynomial representation of sets and additively homomorphic encryption. In [5], a closely related work is proposed by using secure multi-party computation. But the techniques used in their protocols can lead to high communication complexity ( $O(n^2)$ ). In [2], [9], the authors show cryptographic methodologies to compute the over-threshold set-union. Although their protocols have lower communication complexity, their schemes can reveal the multiplicity of data.

Thus, a natural question is: can we design a protocol which has lower communication complexity and can protect data privacy including the multiplicity of data. A positive answer is presented in this paper. We propose a new probabilistic protocol for this problem. Our protocol can ensure that the data aggregator obtains accurate results with high probability (a probability with at least  $1 - 1/n^\epsilon$ , where  $\epsilon \geq 1$ ). Our protocol has nearly linear communication complexity regarding the number of participants while protecting the privacy of participants. We summarize the comparison results in Table 1.

**Table 1.** Complexity Comparisons

Communication	Method	Lower bound	Multiplicity privacy
$O(n^2)$	[10]	$\Omega(n/\log n)^*$ [17]	✓
$O(n^2)$	[5]		✓
$O(n)$	[2]		×
$O(n)$	[9]		×
$O(n(\log n)^2)$	our		✓

\*The lower bound on communication complexity does not consider data privacy and security

Our main contributions in this paper are summarized as follows:

- The proposed protocol can precisely calculate the over-threshold set-union with high probability.
- The proposed protocol not only has  $O(n(\log n)^2)$  communication complexity which nearly matches the lower bound  $\Omega(n/\log n)$  but also protects the data privacy.

The remainder of this paper is organized as follows. We introduce our system model, problem definition, the adversary model, and cryptographic background

in Section 2. In Section 3, we describe our efficient protocol for the over-threshold set-union problem. We show the detailed analysis for our protocol in terms of correctness, security, and complexity in Section 4. Finally, we conclude this paper in Section 5.

## 2 PRELIMINARIES

### 2.1 System Model

In this paper, we adopt the PDA system architecture [2] (see Figure 1), which is composed of three parts: an untrusted data aggregator  $\mathcal{A}_{gg}$ , a proxy  $\mathcal{P}$  and  $n$  participants  $\{p_1, \dots, p_n\}$ . Each participant  $p_i$  holds a private data item  $x_i \in [1, M]$ . The data aggregator has a powerful computational and storage capacity and is responsible for computing some aggregate functions (over-threshold set-union in this paper). Let  $S$  denote the multiset  $\{x_1, \dots, x_n\}$ , and  $[[x_i]]$  denote the number of times of data item  $x_i$  which appears in the multiset  $S$ . For any set  $A$ , let  $x \xleftarrow{U} A$  denote that  $x$  is sampled uniformly at random from  $A$ . We denote by  $|A|$  the number of elements in the set  $A$ . Some symbols used in this paper are listed in Table 2.

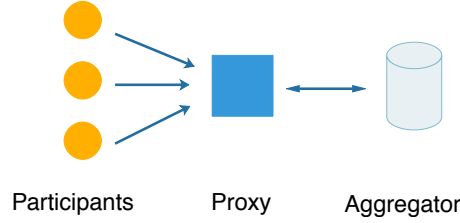


Fig. 1. PDA system architecture

### 2.2 Problem Definition

In the paper, we mainly investigate the over-threshold set-union problem. Informally speaking, this problem is as follows: the aggregator only learns the data items which appear at least  $t$  times in the union of private data items of all participants without obtaining extra knowledge about other data items. The formal definition states as follows.

**Definition 1 (Over-Threshold Set-Union)** *In our system model, each participant  $p_i$  has a data item  $x_i$ , and let  $S = \{x_1, \dots, x_n\}$ . Given a threshold  $t$ , the aggregator wants to compute the set  $\mathcal{U} = \{x'_1, \dots, x'_l\}$  without obtaining extra knowledge, where  $x'_j \in S$ ,  $[[x'_j]] \geq t$  for any  $j \in \{1, \dots, l\}$ .*

**Table 2.** Symbols and Their Definitions

Symbol	Definition
$n$	The number of participants in our model.
$\lambda$	The security parameter.
$[n]$	The set $\{1, \dots, n\}$ .
$t$	Threshold value for over-threshold set union.
$p_i$	The $i$ -th participant, where $i \in [n]$ .
$GL(m, Z_q)$	The general linear group of degree $m$ over $Z_q$ .
$\approx_c$	Computationally indistinguishable.
$x \xleftarrow{U} A$	$x$ is uniformly chosen from set $A$ .
$\mathcal{A}$	Active set about the index of participants.
$\mathcal{U}$	Over-Threshold Set-Union for active participants' data items.
$f_s$	A pseudorandom function indexed by $s$ .
$\Delta_{i,B}(x)$	Lagrange coefficient $= \prod_{j \in B, j \neq i} \frac{x-j}{i-j}$ .

### 2.3 Adversary Model

Like most other related work [10], [2], [9], we consider the honest-but-curious adversaries model. In this model, each party faithfully executes the prescribed actions of the protocol and meanwhile is curious about the private information of other nodes. Moreover, we assume that the participants, proxy and data aggregator do not collude with each other.

### 2.4 Cryptographic background

**Asymmetric Bilinear Pairing Groups** Let us denote the asymmetric bilinear pairing groups by  $(q, \mathbb{G}_1, \mathbb{G}_2, \mathbb{G}_T, e)$  which is produced by  $\mathcal{G}(1^\lambda)$ , where  $q$  is a prime,  $\lambda$  is a security parameter,  $\mathbb{G}_1$ ,  $\mathbb{G}_2$  and  $\mathbb{G}_T$  are cyclic groups of order  $q$ , and  $e: \mathbb{G}_1 \times \mathbb{G}_2 \rightarrow \mathbb{G}_T$ . The bilinear map  $e$  has the following properties:

1. Bilinearity:  $\forall G_1 \in \mathbb{G}_1, \forall G_2 \in \mathbb{G}_2$  and  $\forall a, b \in \mathbb{Z}_q$ ,  $e(aG_1, bG_2) = e(G_1, G_2)^{ab}$ ;
2. Non-degeneracy:  $\forall G_1 \in \mathbb{G}_1$ ,  $e(G_1, G_2) = 1$ , then  $G_2 = 0$ .

**Definition 2** (*Symmetric External Diffie-Hellman: SXDH Assumption [3]*). Given the following distributions  $((q, \mathbb{G}_1, \mathbb{G}_2, \mathbb{G}_T, e) \leftarrow \mathcal{G}(1^\lambda), g_1, g_2, g_1^a, g_1^b, Y_\sigma)$ , where  $Y_0 = g_1^{ab}$ ,  $Y_1 = g_1^{ab+r}$ ,  $a, b, r \xleftarrow{U} \mathbb{Z}_p$ , the advantage  $\mathbf{Adv}_A^{\text{SXDH}}$  that any probabilistic polynomial time (PPT) adversary determines whether  $\sigma = 0$  or 1 is negligible in  $\lambda$ , where  $\mathbf{Adv}_A^{\text{SXDH}} = |\Pr[\mathcal{A}(1^\lambda, Y_0) \rightarrow 1] - \Pr[\mathcal{A}(1^\lambda, Y_1) \rightarrow 1]|$ .

If  $\mathbb{G}_1, \mathbb{G}_2$  reverse the roles, the analogous distributions is also true.

**Definition 3** (*External Decisional Linear Assumption: XDLIN [1]*). Given  $P_a = ((q, \mathbb{G}_1, \mathbb{G}_2, \mathbb{G}_T, e) \leftarrow \mathcal{G}(1^\lambda), \xi G_1, \kappa G_1, \delta \xi G_1, \sigma \kappa G_1, \xi G_2, \kappa G_2, \delta \xi G_2, \sigma \kappa G_2, Y_a)$  for  $b \in \{1, 2\}$ , where  $Y_0 = (\delta + \sigma)G_b$ ,  $Y_1 = (\delta + \sigma + \rho)G_b$ ,  $\delta, \kappa, \sigma, \rho \xleftarrow{U} \mathbb{Z}_p$ , the advantage  $\mathbf{Adv}_A^{\text{XDLIN}}$  that any PPT adversary determines whether  $a = 0$  or 1 is negligible in  $\lambda$ , where  $\mathbf{Adv}_A^{\text{XDLIN}} = |\Pr[\mathcal{A}(1^\lambda, Y_0) \rightarrow 1] - \Pr[\mathcal{A}(1^\lambda, Y_1) \rightarrow 1]|$ .

**Definition 4 (Dual Pairing Vector Spaces [14] [16])** *Dual Pairing Vector Spaces are defined by the tuple  $(q, \mathbb{V}, \mathbb{V}^*, \mathbb{G}_T, \mathbb{A}, \mathbb{A}^*)$  and can be constructed by  $\mathcal{G}(1^\lambda)$ .  $\mathbb{V} = \mathbb{G}_1^m$  and  $\mathbb{V}^* = \mathbb{G}_2^m$  are  $n$  dimensional vector spaces and  $\mathbb{V} = (\mathbf{a}_1, \dots, \mathbf{a}_m)$  and  $\mathbb{V}^* = (\mathbf{a}_1^*, \dots, \mathbf{a}_m^*)$  are canonical bases of  $\mathbb{V}$  and  $\mathbb{V}^*$ , respectively, where  $\mathbf{a}_i = (0^{i-1}, G_1, 0^{m-i})$ ,  $\mathbf{a}_i^* = (0^{i-1}, G_2, 0^{m-i})$ , where  $0^j$  represents a line of  $j$  zeros, e.g.,  $(0^4, 0) = (0, 0, 0, 0, 0)$ .  $\tilde{e} : \mathbb{V} \times \mathbb{V}^* \rightarrow \mathbb{G}_T$  is pairing and is defined by  $\tilde{e} = (\mathbf{x}, \mathbf{y}) = \prod_{i=1}^m e(X_i, Y_i) \in \mathbb{G}_T$ , where  $\mathbf{x} = (X_1, \dots, X_m) \in \mathbb{V}$  and  $\mathbf{y} = (Y_1, \dots, Y_m) \in \mathbb{V}^*$ .*

Then, given  $(q, \mathbb{V}, \mathbb{V}^*, \mathbb{G}_T, \mathbb{A}, \mathbb{A}^*, \tilde{e})$ , we present random dual orthonormal bases  $\mathcal{G}_{ob}(1^\lambda, m)$  as follows:

$\psi \xleftarrow{U} Z_p, \mathbf{B} = (\chi_{i,j}) \leftarrow GL(m, Z_q), (\phi_{i,j}) \leftarrow \psi(\mathbf{B}^T)^{-1}, \mathbf{b}_i = \sum_{j=1}^m \chi_{i,j} \mathbf{a}_j, \mathbf{b}_i^* = \sum_{j=1}^m \phi_{i,j} \mathbf{a}_j^*, (i = 1, \dots, m), \mathbb{B} = (\mathbf{b}_1, \dots, \mathbf{b}_m), \mathbb{B}^* = (\mathbf{b}_1^*, \dots, \mathbf{b}_m^*), g_T = e(G_1, G_2)^\psi$ ,  $\text{output}(\mathbb{B}, \mathbb{B}^*)$ , where  $GL(m, Z_q)$  denotes the general linear group of degree  $m$  over  $Z_p$ .

For  $\mathbf{x} = (X_1, \dots, X_m)^T$  and  $\mathbb{B} = (\mathbf{b}_1, \dots, \mathbf{b}_m)$ , we use  $(\mathbf{x})_{\mathbb{B}}$  to denote  $\sum_{i=1}^m X_i \mathbf{b}_i$ . So we have  $\tilde{e}((\mathbf{x})_{\mathbb{A}}, (\mathbf{y})_{\mathbb{A}^*}) = \prod_{i=1}^m e(X_i G_1, Y_i G_2) = e(G_1, G_2)^{\sum_{i=1}^m X_i Y_i}$  and  $\tilde{e}((\mathbf{x})_{\mathbb{B}}, (\mathbf{y})_{\mathbb{B}^*}) = \tilde{e}((\mathbf{B}\mathbf{x})_{\mathbb{A}}, (\psi(\mathbf{B}^T)^{-1}\mathbf{y})_{\mathbb{A}^*}) = e(G_1, G_2)^{\psi \mathbf{B}\mathbf{x} \cdot \mathbf{B}^{-1}\mathbf{y}} = g_T^{\sum_{i=1}^m X_i Y_i}$ .

**Threshold Secret Sharing** In this paper, we use an  $(n, t)$ -threshold secret sharing scheme, i.e., Shamir's secret sharing technique [15]. The secret sharing scheme is that a secret  $s$  is split and distributed to  $n$  participants, and only  $t$  or more than  $t$  participants which work together can reconstruct the secret  $s$ . For the secret  $s$ , it constructs polynomial  $h(x) = s + s_1x + s_2x^2 + \dots + s_{t-1}x^{t-1}$ , where  $s_1, \dots, s_{t-1} \xleftarrow{U} Z_p$ . Then, it randomly chooses  $n$  distinct elements  $\alpha_1, \dots, \alpha_n$  and computes  $h(\alpha_1), \dots, h(\alpha_n)$  from  $Z_p$  and sends  $(\alpha_i, h(\alpha_i))$  to participant  $p_i$  for any  $i \in [n]$ . For any set  $B = \{\alpha_{i_1}, \dots, \alpha_{i_t}\}$  of size  $t$ , the secret  $s$  can be recovered as followings:  $s = \sum_{j=1}^t \Delta_{\alpha_{i_j}, B}(0) h(\alpha_{i_j})$ , where  $\Delta_{\alpha_{i_j}, B}(0) = \prod_{r \in B, r \neq \alpha_{i_j}} \frac{-r}{\alpha_{i_j} - r}$ .

Moreover, we introduce the distributed ElGamal cryptosystem [7] in a threshold manner, which will be adopted in our paper. A simplified description of the cryptosystem between the proxy  $\mathcal{P}$  and data aggregator  $\mathcal{A}_{gg}$  is described as follows. Let  $\mathbb{G}$  be a multiplicative cyclic group of prime order  $q$  and  $g$  is one of its generator.  $\mathcal{P}$  and  $\mathcal{A}_{gg}$  choose private keys  $sk_1$  and  $sk_2$  from  $Z_p$ , respectively. Then, they publish  $y_1 = g^{sk_1}$  and  $y_2 = g^{sk_2}$ , and compute the common public key  $pu = y_1 y_2 = g^{sk_1 + sk_2}$ . In order to encrypt a message  $m$ , any participant chooses  $r \xleftarrow{U} Z_p$  and computes a ciphertext  $\mathbf{E}_{pu}[m] = (A, B) = (g^r, my^r)$ . For decryption, on inputs a ciphertext  $(A, B)$ ,  $\mathcal{P}$  sends  $A_1 = A^{sk_1}$  to  $\mathcal{A}_{gg}$  and  $\mathcal{A}_{gg}$  can obtain the message  $m$  by computing  $B / (A_1 A^{sk_2})$  (it is not difficult to see that  $B / (A_1 A^{sk_2}) = my^r / (g^{rsk_1} g^{rsk_2}) = my^r / g^{(sk_1 + sk_2)r} = my^r / y^r = m$ ). The distributed ElGamal cryptosystem has one important property: re-encryption. That is, given  $\mathbf{E}_{pu}[m] = (A, B)$ , the proxy can compute another encryption of  $m$ :  $(Ag^{r'}, By^{r'})$ , which is still denoted as  $\mathbf{E}_{pu}[m]$  for convenience.

### 3 Efficient Protocol for Threshold-Over Set-Union

#### 3.1 Overview

In the previous schemes [10], [5], [2], [9], they process the data items directly. And yet we use the binary bits of data items to compute the over-threshold set-union. The main idea for our protocol is that, only  $t$  or more than  $t$  ciphertexts which are generated by using the same data item can be decrypted. To achieve the design goal and reduce the communication complexity, we propose a novel protocol based on the threshold secret sharing and dual pairing vector spaces. Our protocol is composed of the four phases as follows: the **System Setup**, **Ciphertext Generation**, **Blind Ciphertext** and **Ciphertext Aggregation**.

#### 3.2 Basic Protocol

We first introduce a basic protocol which can exactly compute the over-threshold set-union with high probability.

**System Setup** A trusted authority produces some private keys and distributes to the participants, and the data aggregator and proxy generate their own secret keys. Specifically, it selects  $(q, \mathbb{G}_1, \mathbb{G}_2, \mathbb{G}_T, e) \leftarrow \mathcal{G}(1^\lambda)$  and generates  $(q, \mathbb{V}, \mathbb{V}^*, \mathbb{G}_T, \mathbb{A}, \mathbb{A}^*, \tilde{e}), \{(\mathbb{B}_j, \mathbb{B}_j^*), (\mathbb{C}_j, \mathbb{C}_j^*)\}_{j=1}^3$  from  $\mathcal{G}_{ob}(1^\lambda, m)$ , where  $\mathbb{B}_j = (\mathbf{b}_{j,1}, \dots, \mathbf{b}_{j,9}), \mathbb{B}_j^* = (\mathbf{b}_{j,1}^*, \dots, \mathbf{b}_{j,9}^*)$  for  $j = 1, 2$ ,  $\mathbb{B}_3 = (\mathbf{b}_{3,1}, \dots, \mathbf{b}_{3,6l+11}), \mathbb{B}_3^* = (\mathbf{b}_{3,1}^*, \dots, \mathbf{b}_{3,6l+11}^*), \mathbb{C}_j = (\mathbf{c}_{j,1}, \dots, \mathbf{c}_{j,6l+9}), \mathbb{C}_j^* = (\mathbf{c}_{j,1}^*, \dots, \mathbf{c}_{j,6l+9}^*)$  for  $j = 1, 2$ ,  $\mathbb{C}_3 = (\mathbf{c}_{3,1}, \dots, \mathbf{c}_{3,6l+11}), \mathbb{C}_3^* = (\mathbf{c}_{3,1}^*, \dots, \mathbf{c}_{3,6l+11}^*)$ . Then it computes  $\{(\alpha_{i,j}, h_{i,j}), (\alpha_{i,j}, \bar{h}_{i,j})\}$  for secrets  $\{\beta_j, \gamma_j\}$  by using the Shamir's protocol and samples  $\{a_j, b_j\}$  from  $Z_q$  (to be used in pseudo-random function (PRF)), where  $i \in [n], j \in [2] = \{1, 2\}$ . It sends  $(q, \mathbb{G}_1, \mathbb{G}_2, \mathbb{G}_T, e), \{(\mathbb{B}_j, \mathbb{B}_j^*), (\mathbb{C}_j, \mathbb{C}_j^*)\}_{j=1}^3, \{a_1, b_1\}, \{(\alpha_{i,j}, h_{i,j}), (\alpha_{i,j}, \bar{h}_{i,j})\}_{j=1}^2$  to participant  $p_i$  for any  $i \in [n]$  and  $(q, \mathbb{G}_1, \mathbb{G}_2, \mathbb{G}_T, e)$  to the data aggregator and proxy.

The data aggregator together with proxy compute secret key  $sk_1$  and  $sk_2$ , respectively, and they publish common public key  $pu$  by running the distributed ElGamal cryptosystem.

**Ciphertext Generation** In this phase, each participant  $p_i$  encrypts its own data item  $x_i$  by using private keys.

Specifically, participant  $p_i$  first converts  $x_i$  into binary vector  $\mathbf{v}_i = (x_{i,l-1}, \dots, x_{i,0})$ , where  $l = \lceil \log M \rceil$ . Let  $\mathbf{X}_{i,1}$  and  $\mathbf{X}_{i,2}$  denote the vectors  $(\mathbf{v}_i, \mathbf{1}, \mathbf{v}_i)$  and  $(\mathbf{1}, \mathbf{v}_i, -2\mathbf{v}_i)$ , respectively.  $f_{a_1}(\cdot) (f_{b_1}(\cdot))$  is a pseudorandom function of the PRF family  $\mathcal{F}_\lambda = \{f_s : \{0, 1\}^\lambda \rightarrow \{0, 1\}^\lambda\}_{s \in \{0, 1\}^\lambda}$  that uses  $a_1$  ( $b_1$ ) as parameter [6]. Then it computes the ciphertexts  $\bar{C}_i = \{C_{i,1}, \dots, C_{i,6}, C_i\}$  as follows:

$$\begin{aligned}
C_{i,1} &= (r_{i,1}, h_{i,1}, 0^2, r_{i,2}, r_{i,3}, 0, 0, 0)_{\mathbb{B}_1}, C_{i,2} = (r_{i,4}, r_{i,5}, 0^2, 0, 0, r_{i,6}, r_{i,7}, 0)_{\mathbb{B}_1^*}, \\
C_{i,3} &= (r_{i,4}, \bar{h}_{i,1}, 0^2, \hat{r}_{i,1}, \hat{r}_{i,2}, 0, 0, 0)_{\mathbb{B}_2}, C_{i,4} = (r_{i,1}, \hat{r}_{i,3}, 0^2, 0, 0, \hat{r}_{i,4}, \hat{r}_{i,5}, 0)_{\mathbb{B}_2^*}, \\
C_{i,5} &= (\bar{r}_{i,1} f_{a_1}(r) \cdot \mathbf{X}_{i,1}, r_{i,5}, \hat{r}_{i,3}, 1, 0^{3l+3}, \bar{r}_{i,2}, \bar{r}_{i,3}, 0, 0, 0)_{\mathbb{B}_3}, \\
C_{i,6} &= (f_{b_1}(r) \cdot \mathbf{X}_{i,2}, h_{i,1}, \bar{h}_{i,1}, \bar{r}_{i,4}, 0^{3l+3}, 0, 0, \bar{r}_{i,5}, \bar{r}_{i,6}, 0)_{\mathbb{B}_3^*}, \\
C_i &= e(G_1, G_2)^{\beta_1 r_{i,5} + \gamma_1 \hat{r}_{i,3} + \bar{r}_{i,4}},
\end{aligned}$$

where  $\{r_{i,j}\}_{j=1}^7, \{\hat{r}_{i,j}\}_{j=1}^5, \{\bar{r}_{i,j}\}_{j=2}^6 \xleftarrow{U} Z_q, \bar{r}_{i,1} \xleftarrow{U} Z_{q'} (q' < q/tl), \mathbf{1} = (1, 1, \dots, 1)$ , and  $r$  is a nonce about the message, and  $f_{a_1}(r)f_{b_1}(r) \neq 0$ . Similarly, it computes the following ciphertexts  $\bar{D}_i = \{\{D_{i,j}\}_{j=1}^6, D_i\}$ :

$$\begin{aligned}
D_{i,1} &= (s_{i,1} \cdot \mathbf{X}_{i,1}, s_{i,2}, h_{i,2}, 0^{3l+2}, s_{i,3}, s_{i,4}, 0, 0, 0)_{\mathbb{C}_1}, \\
D_{i,2} &= (s_{i,5} \cdot \mathbf{X}_{i,2}, s_{i,6}, s_{i,7}, 0^{3l+2}, 0, 0, s_{i,8}, s_{i,9}, 0)_{\mathbb{C}_1^*}, \\
D_{i,3} &= (\hat{s}_{i,1} \cdot \mathbf{X}_{i,1}, s_{i,6}, \bar{h}_{i,2}, 0^{3l+2}, \hat{s}_{i,2}, \hat{s}_{i,3}, 0, 0, 0)_{\mathbb{C}_2}, \\
D_{i,4} &= (\hat{s}_{i,4} \cdot \mathbf{X}_{i,2}, s_{i,2}, \hat{s}_{i,5}, 0^{3l+2}, 0, 0, \hat{s}_{i,6}, \hat{s}_{i,7}, 0)_{\mathbb{C}_2^*}, \\
D_{i,5} &= (\bar{s}_{i,1} \cdot \mathbf{X}_{i,1}, s_{i,7}, \hat{s}_{i,5}, 1, 0^{3l+3}, \bar{s}_{i,2}, \bar{s}_{i,3}, 0, 0, 0)_{\mathbb{C}_3}, \\
D_{i,6} &= (\bar{s}_{i,4} \cdot \mathbf{X}_{i,2}, h_{i,2}, \bar{h}_{i,2}, \bar{s}_{i,5}, 0^{3l+3}, 0, 0, \bar{s}_{i,6}, \bar{s}_{i,7}, 0)_{\mathbb{C}_3^*}, \\
D_i &= e(G_1, G_2)^{x_i + \beta_2 s_{i,7} + \gamma_2 \hat{s}_{i,5} + \bar{s}_{i,5}},
\end{aligned}$$

where  $\{s_{i,j}\}_{j=1}^9, \{\hat{s}_{i,j}, \bar{s}_{i,j}\}_{j=1}^7 \xleftarrow{U} Z_q$ .

Then,  $p_i$  handles  $\bar{C}_i$  and  $\bar{D}_i$  by using  $pu$  according to the distributed ElGamal cryptosystem, and obtains  $E_i = (\mathbf{E}_{pu}[\bar{C}_i], \{\mathbf{E}_{pu}[C_{i,j}]\}_{j=1}^6, \mathbf{E}_{pu}[\alpha_{i,1}])$  and  $E'_i = (\mathbf{E}_{pu}[\bar{D}_i], \{\mathbf{E}_{pu}[D_{i,j}]\}_{j=1}^6, \mathbf{E}_{pu}[\alpha_{i,2}])$ .

Finally,  $p_i$  sends  $\{E_i, E'_i\}$  to the proxy.

**Blind Ciphertext** In this phase, the proxy needs to blind the ciphertexts such that the data aggregator does not learn knowledge about the source of ciphertexts. In fact, communication is very expensive and the participants drop out at any time in many important applications. Thus, denote the index  $i$  of the participants  $p_i$  that the proxy receives messages by active set  $\mathcal{A}$ . After receiving messages from  $\mathcal{A}$ , the proxy performs the blind operation as follows.

It first performs two shuffle operations for the ciphertext  $\{E_i\}_{i \in \mathcal{A}}$  and  $\{E'_i\}_{i \in \mathcal{A}}$ . Specifically, the proxy picks two random permutations  $\pi$  and  $\pi'$  from the set  $\mathcal{A}$ , and re-encrypts  $E = \{E_{\pi(i)}\}_{i \in \mathcal{A}}$  and  $E' = \{E'_{\pi'(i)}\}_{i \in \mathcal{A}}$  by using the re-encryption property of distributed ElGamal cryptosystem.

Finally, it sends  $\{E, E'\}$  to the data aggregator.

**Ciphertext Aggregation** In this phase,  $\mathcal{A}_{gg}$  computes the over-threshold set-union by using decryption operation. The detailed decryption operation is summarized in Algorithm 1.

---

**Algorithm 1: Ciphertext Aggregation:  $\mathcal{A}_{gg}$** 


---

**Input:**  $\{E_i, \bar{E}_i\}_{i \in \mathcal{A}}, \mathcal{U} \leftarrow \phi$   
**Output:** Over-Threshold Set-Union  $\mathcal{U}$

- 1 with the help of the proxy,  $\mathcal{A}_{gg}$  decrypts  $\{E_i, \bar{E}_i\}_{i \in \mathcal{A}}$  based on the distributed ElGamal cryptosystem and obtains  $\{\bar{C}_j, \alpha_{j,1}\}_{j \in \mathcal{A}}, \{\bar{D}_j, \alpha_{j,2}\}_{j \in \mathcal{A}};$
- 2 **if**  $|\mathcal{A}| < t$  **then**
- 3      $\mathcal{A}_{gg}$  outputs an empty set  $\phi;$
- 4 **else**
- 5      $\mathcal{A}_{gg}$  computes a family of sets  $\mathcal{U}_1$  whose elements are all subsets of size  $t$  of  $\mathcal{A}$ , where  $|\mathcal{U}_1| = \binom{|\mathcal{A}|}{t};$
- 6     **while**  $|\mathcal{U}_1| \leq 0$  **do**
- 7          $\mathcal{A}_{gg}$  randomly chooses an element  $\mathcal{E} = \{e_1, \dots, e_t\}$  from  $\mathcal{U}_1;$
- 8          $\mathcal{U}_1 \leftarrow \mathcal{U}_1 - \{\mathcal{E}\}, B \leftarrow \{\alpha_{i,1} | i \in \mathcal{E}\}, B' \leftarrow \{\alpha_{i,2} | i \in \mathcal{E}\};$
- 9          $\mathcal{A}_{gg}$  computes  $H = \prod_{j \in \mathcal{E}} C_j, F_{i,1} = \prod_{j \in \mathcal{E}} e(C_{i,1}, C_{j,2}),$   
 $F_{i,2} = \prod_{j \in \mathcal{E}} e(C_{i,4}, C_{j,3}), \hat{F}_{i,2} = \prod_{j \in \mathcal{E}} e(C_{i,3}, C_{j,4}),$   
 $\hat{F}_{i,2} = \prod_{j \in \mathcal{E}} e(C_{i,2}, C_{j,1}), \bar{F}_i = \prod_{j \in \mathcal{E}} e(C_{i,5}, C_{j,6});$
- 10         **if**  $\prod_{i \in \mathcal{E}} \left( \frac{F_{i,1} \hat{F}_{i,1} \bar{F}_i}{F_{i,2} \hat{F}_{i,2}} \right)^{\Delta_{\alpha_{i,1}, B}} == H$  **then**
- 11              $\mathcal{A}_{gg}$  computea  $K = \prod_{j \in \mathcal{E}} D_j, J_{i,1} = \prod_{j \in \mathcal{E}} e(D_{i,1}, D_{j,2}),$   
 $J_{i,2} = \prod_{j \in \mathcal{E}} e(D_{i,4}, D_{j,3}), \hat{J}_{i,2} = \prod_{j \in \mathcal{E}} e(D_{i,3}, D_{j,4}),$   
 $\hat{J}_{i,2} = \prod_{j \in \mathcal{E}} e(D_{i,2}, D_{j,1}), \bar{J}_i = \prod_{j \in \mathcal{E}} e(D_{i,5}, D_{j,6});$
- 12              $\hat{U} = \left( \prod_{i \in \mathcal{E}} K_i \left( \frac{J_{i,2} \hat{J}_{i,2}}{J_{i,1} \bar{J}_i} \right)^{\Delta_{\alpha_{i,2}, B'(0)}} \right)^{\frac{1}{t}}, U = \log \hat{U};$
- 13              $\mathcal{U} \leftarrow \mathcal{U} \cup U;$
- 14 **return**  $\mathcal{U}$

---



With the help of the proxy,  $\mathcal{A}_{gg}$  decrypts the ciphertexts  $E$  and  $E'$  and obtains  $\{\overline{C}_j, \alpha_{j,1}\}_{j \in \mathcal{A}}$ ,  $\{\overline{D}_j, \alpha_{j,2}\}_{j \in \mathcal{A}}$  (line 1). Next, if  $|\mathcal{A}| < t$ , the participants' data items do not contain the over-threshold set-union and the data aggregator outputs an empty set  $\phi$  (lines 2-3). Otherwise, they may contain the over-threshold set-union. Any  $t$  participants' data items could be the same.  $\mathcal{A}_{gg}$  needs to traverse all possible combination of  $t$  participants' ciphertexts. Thus, it computes a family of sets  $\mathcal{U}_1$  whose elements are all subsets of size  $t$  of  $\mathcal{A}$  (line 5). Obviously,  $|\mathcal{U}_1| = \binom{|\mathcal{A}|}{t}$ . Then it randomly chooses an element  $\mathcal{E} = \{e_1, \dots, e_t\}$  from  $\mathcal{A}$  and removes  $\mathcal{E}$  from  $\mathcal{A}$  and computes  $B$ ,  $B'$  (line 8). Next,  $\mathcal{A}_{gg}$  computes  $H, F_{i,1}, F_{i,2}, \hat{F}_{i,1}, \hat{F}_{i,2}$  and  $\overline{F}_i$  (lines 9). Then, it calculates  $\prod_{i \in \mathcal{E}} \left( \frac{F_{i,1} \hat{F}_{i,1} \overline{F}_i}{F_{i,2} \hat{F}_{i,2}} \right)^{\Delta_{\alpha_{i,1}, B}(0)}$  and  $H$ , and decides whether  $\prod_{i \in \mathcal{E}} \left( \frac{F_{i,1} \hat{F}_{i,1} \overline{F}_i}{F_{i,2} \hat{F}_{i,2}} \right)^{\Delta_{\alpha_{i,1}, B}(0)} = H$  (line 10).

If the above equation holds, the data aggregator continues to compute  $K, J_{i,1}, J_{i,2}, \hat{J}_{i,1}, \hat{J}_{i,2}, \overline{J}_i$  and  $\hat{U}$  (lines 11-12). To recover  $U$ , it needs to compute the discrete logarithm of  $\hat{U}$  to the base  $e(G_1, G_2)$ . By using the Pollard's lambda algorithm [13], this takes expected time  $O(\sqrt{M})$  since  $x_i \in [1, M]$ . Note that if the above equation (line 10) holds, then it means the  $t$  ciphertexts contain the same data item which belongs to over-threshold set-union (see Theorem 1). We can obtain the data item by using the above operations.

Finally, to obtain over-threshold set-union, the data aggregator requires at most  $\binom{n}{t}$  times of the above computation (lines 6-13), because  $\mathcal{U}_1$  contains at most  $\binom{n}{t}$  elements when  $|\mathcal{A}| = n$ .

### 3.3 Further Improvement

In fact, each  $p_i$  generates a new private data item in a certain time interval. If participants use the basic protocol to encrypt their data items in each time interval, the data aggregator may learn more information than the over-threshold set-union. For example, suppose that  $p_1$  produces the same data item in  $t$  time intervals and encrypts it by using the basic protocol, and if the data item does not belong to over-threshold set-union in the  $t$  time intervals, the data aggregator can obtain the data item by using the decryption operation. Thus, the basic protocol cannot be applied to this scenario.

To resolve this drawback, we propose our improved algorithm in this section. The main idea is that the ciphertexts can be decrypted only if they are at the same time interval. Specifically, since ciphertexts  $C_{i,j}$  and  $D_{i,j}$  are computed by using some vectors and random dual orthonormal bases in each time interval, we can use different vectors in different time interval. An efficient method is that we can rearrange the order of the elements in the vectors. For example,  $D_{i,1}$  is computed by using vector  $(s_{i,1} \cdot \mathbf{X}_{i,1}, s_{i,2}, h_{i,2}, 0^{3l+2}, s_{i,3}, s_{i,4}, 0, 0, 0)$ , and we use a new vector whose elements are randomly arranged to compute ciphertexts in each time interval. Although the scheme is feasible, there is a disadvantage to this approach. Because the number of elements in vectors is limited. Thus, the number of random permutation is limited. To overcome the weakness, we can

add some elements to vectors. In the following, we will illustrate our improved algorithm in detail.

We again run the **System Setup** phase for the trusted authority, but we make the following modifications.

- Add: samples  $c_{i,j}, d_i \leftarrow \{0,1\}^\lambda$  ( $i = 1, \dots, 6, j = 1, \dots, 4$ ) and sends to all participants;
- Replace  $\{(\mathbb{B}_i, \mathbb{B}_i^*), (\mathbb{C}_i, \mathbb{C}_i^*)\}_{i=1}^3$  by  $\{(\overline{\mathbb{B}}_i, \overline{\mathbb{B}}_i^*), (\overline{\mathbb{C}}_i, \overline{\mathbb{C}}_i^*)\}_{i=1}^3$  : Samples  $\{(\overline{\mathbb{B}}_i, \overline{\mathbb{B}}_i^*), (\overline{\mathbb{C}}_i, \overline{\mathbb{C}}_i^*)\}_{i=1}^3 \leftarrow \mathcal{G}_{ob}(1^\lambda, m)$ , where  $m = 6l + 11 + c_\eta \log n + c_u$ ;

To rearrange the order of the elements in the vectors in different time intervals, we use a pseudorandom permutation and the Luby-Rackoff construction [12]. Let  $\overline{f}_{s_1, s_2, s_3, s_4}()$  be a pseudorandom permutation indexed by  $s_1, s_2, s_3, s_4$  (because it uses four pseudorandom functions), where  $\overline{f}_{s_1, s_2, s_3, s_4}() : \{0,1\}^{6l+11+c_\eta \log n + c_u} \times \{0,1\}^\lambda \rightarrow \{0,1\}^{6l+11+c_\eta \log n + c_u}$ , and  $c_\eta$  and  $c_u$  are large constants. For a vector  $(x_1, x_2, x_3)$ , we use positive integers 1, 2, 3 to label the positions of the elements  $x_1, x_2, x_3$  in the vector from left to right and let  $R(x_1, x_2, x_3)$  denote them, i.e.,  $R(x_1, x_2, x_3) = (1, 2, 3)$ . Let  $\overline{f}_{s_1, s_2, s_3, s_4}(R(x_1, x_2, x_3))$  denote  $\overline{f}_{s_1, s_2, s_3, s_4}(R(x_1, x_2, x_3)) = \overline{f}_{s_1, s_2, s_3, s_4}(1, 2, 3) = (\overline{f}_{s_1, s_2, s_3, s_4}(1), \overline{f}_{s_1, s_2, s_3, s_4}(2), \overline{f}_{s_1, s_2, s_3, s_4}(3))) = (y_1, y_2, y_3)$ . Let  $\overline{R}(y_1, y_2, y_3)$  denote a new vector which consists of  $x_1, x_2, x_3$  and rearranges the positions of  $x_1, x_2, x_3$  according to  $y_1, y_2, y_3$ , respectively. For example, when  $(y_1, y_2, y_3) = (3, 2, 1)$ , we have  $\overline{R}(y_1, y_2, y_3) = (x_3, x_2, x_1)$ . For convenience, we use  $\langle k_i(r) \rangle$  to denote  $f_{d_{i,1}}(r), f_{d_{i,2}}(r), f_{d_{i,3}}(r), f_{d_{i,4}}(r)$  for  $i = 1, \dots, 6$ .

We again run the **Ciphertext Generation** phase for  $p_i$  ( $i \in [n]$ ), but we make the following modifications.

- Replace  $\{\{C_{i,j}\}_{j=1}^6, C_i\}$  by  $\{\{\tilde{C}_{i,j}\}_{j=1}^6, \tilde{C}_i\}$  :  

$$\tilde{C}_{i,1} = \overline{R}\left(\overline{f}_{\langle k_1(r) \rangle}(R(r_{i,1}, h_{i,1}, 0^2, r_{i,2}, r_{i,3}, 0, 0, 0, f_{c_1}(1), \dots, f_{c_1}(t'))))\right)_{\overline{\mathbb{B}}_1},$$

$$\tilde{C}_{i,2} = \overline{R}\left(\overline{f}_{\langle k_1(r) \rangle}(R(r_{i,4}, r_{i,5}, 0^2, 0, 0, r_{i,6}, r_{i,7}, 0, u_{i,1}, \dots, u_{i,t'}))\right)_{\overline{\mathbb{B}}_1^*},$$

$$\tilde{C}_{i,3} = \overline{R}\left(\overline{f}_{\langle k_2(r) \rangle}(R(r_{i,4}, \bar{h}_{i,1}, 0^2, \hat{r}_{i,1}, \hat{r}_{i,2}, 0, 0, 0, f_{c_2}(1), \dots, f_{c_2}(t'))))\right)_{\overline{\mathbb{B}}_2},$$

$$\tilde{C}_{i,4} = \overline{R}\left(\overline{f}_{\langle k_2(r) \rangle}(R(r_{i,1}, \hat{r}_{i,3}, 0^2, 0, 0, \hat{r}_{i,4}, \hat{r}_{i,5}, 0, \hat{u}_{i,1}, \dots, \hat{u}_{i,t'}))\right)_{\overline{\mathbb{B}}_2^*},$$

$$\tilde{C}_{i,5} = \overline{R}\left(\overline{f}_{\langle k_3(r) \rangle}(R(\bar{r}_{i,1} f_{a_1}(r) \cdot \mathbf{X}_{i,1}, r_{i,5}, \hat{r}_{i,3}, 1, 0^{3l+3}, \bar{r}_{i,2}, \bar{r}_{i,3}, 0, 0, 0, f_{c_3}(1), \dots, f_{c_3}(t''))\right)_{\overline{\mathbb{B}}_3},$$

$$\tilde{C}_{i,6} = \overline{R}\left(\overline{f}_{\langle k_3(r) \rangle}(R(f_{b_1}(r) \cdot \mathbf{X}_{i,2}, h_{i,1}, \bar{h}_{i,1}, \bar{r}_{i,4}, 0^{3l+3}, 0, 0, \bar{r}_{i,5}, \bar{r}_{i,6}, 0, \bar{u}_{i,1}, \dots, \bar{u}_{i,t''}))\right)_{\overline{\mathbb{B}}_3^*},$$

$$\tilde{C}_i = e(G_1, G_2)^{\beta_1 r_{i,5} + \gamma_1 \hat{r}_{i,3} + \bar{r}_{i,4}}$$

where  $t' = 6l + 2 + c_\eta \log n + c_u$ ,  $t'' = c_\eta \log n + c_u$ , and  $u_{i,j}, \hat{u}_{i,j}, \bar{u}_{i,j'} \leftarrow Z_p$  for  $j = 1, \dots, t' - 1$  and  $j' = 1, \dots, t'' - 1$ , and  $\sum_{j=1}^{t'} f_{c_1}(j) u_{i,j} = 0$ ,  $\sum_{j=1}^{t'} f_{c_2}(j) \hat{u}_{i,j} = 0$ ,  $\sum_{j=1}^{t''} f_{c_3}(j) \bar{u}_{i,j} = 0$

- Replace  $\{\{D_{i,j}\}_{j=1}^6, D_i\}$  by  $\{\{\tilde{D}_{i,j}\}_{j=1}^6, \tilde{D}_i\}$ :  

$$\tilde{D}_{i,1} = \overline{R}\left(\overline{f}_{\langle k_4(r) \rangle}(R(s_{i,1} \cdot \mathbf{X}_{i,1}, s_{i,2}, h_{i,2}, 0^{3l+2}, s_{i,3}, s_{i,4}, 0, 0, 0, f_{c_4}(1), \dots, f_{c_4}(t'))\right)_{\overline{\mathbb{C}}_1},$$

$$\begin{aligned}
\tilde{D}_{i,2} &= \overline{R}(\overline{f}_{<k_4(r)>}(R(s_{i,5} \cdot \mathbf{X}_{i,2}, s_{i,6}, s_{i,7}, 0^{3l+2}, 0, 0, s_{i,8}, s_{i,9}, 0, j_{i,1}, \dots, j_{i,\hat{t}'})))_{\overline{\mathbb{C}}_1^*}, \\
\tilde{D}_{i,3} &= \overline{R}(\overline{f}_{<k_5(r)>}(R(\hat{s}_{i,1} \cdot \mathbf{X}_{i,1}, s_{i,6}, \bar{h}_{i,2}, 0^{3l+2}, \hat{s}_{i,2}, \hat{s}_{i,3}, 0, 0, 0, f_{c_5}(1), \dots, f_{c_5}(\hat{t}'))))_{\overline{\mathbb{C}}_2}, \\
\tilde{D}_{i,4} &= \overline{R}(\overline{f}_{<k_5(r)>}(R(\hat{s}_{i,4} \cdot \mathbf{X}_{i,2}, s_{i,2}, \hat{s}_{i,5}, 0^{3l+2}, 0, 0, \hat{s}_{i,6}, \hat{s}_{i,7}, 0, \hat{j}_{i,1}, \dots, \hat{j}_{i,\hat{t}'})))_{\overline{\mathbb{C}}_2^*}, \\
\tilde{D}_{i,5} &= \overline{R}(\overline{f}_{<k_6(r)>}(R(\bar{s}_{i,1} \cdot \mathbf{X}_{i,1}, s_{i,7}, \hat{s}_{i,5}, 1, 0^{3l+3}, \bar{s}_{i,2}, \bar{s}_{i,3}, 0, 0, 0, f_{c_6}(1), \dots, f_{c_6}(\hat{t}'))))_{\overline{\mathbb{C}}_3}, \\
\tilde{D}_{i,6} &= \overline{R}(\overline{f}_{<k_6(r)>}(R(\bar{s}_{i,4} \cdot \mathbf{X}_{i,2}, h_{i,2}, \bar{h}_{i,2}, \bar{s}_{i,5}, 0^{3l+3}, 0, 0, \bar{s}_{i,6}, \bar{s}_{i,7}, 0, \bar{j}_{i,1}, \dots, \bar{j}_{i,\hat{t}''})))_{\overline{\mathbb{C}}_3^*}, \\
\tilde{D}_i &= e(G_1, G_2)^{x_i + \beta_2 s_{i,7} + \gamma_2 \hat{s}_{i,5} + \bar{s}_{i,5}}, \\
\text{where } \hat{t}' &= 2 + c_\eta \log n + c_u \text{ and } \hat{t}'' = c_\eta \log n + c_u, \text{ and } j_{i,k}, \hat{j}_{i,k}, \bar{j}_{i,k'} \leftarrow \\
Z_p &\text{ for } k = 1, \dots, \hat{t}' - 1 \text{ and } k' = 1, \dots, \hat{t}'' - 1, \text{ and } \sum_{k=1}^{\hat{t}'} f_{c_1}(k) j_{i,k} = 0, \\
\sum_{k=1}^{\hat{t}'} f_{c_2}(k) \hat{j}_{i,k} &= 0, \sum_{k=1}^{\hat{t}''} f_{c_3}(k) \bar{j}_{i,k} = 0
\end{aligned}$$

Note that we add some elements to the vectors used to compute ciphertexts such that the vectors have the same number of elements. Moreover, the added elements do not affect the result of executed protocol. For  $O(\log n)$  elements, we have that the number of different permutation is  $O((\log n)!)$ . By using Stirling's approximation, we have  $O((\log n)!) \sim O(\sqrt{2\pi \log n} (\frac{\log n}{e})^{\log n})$ . On the other hand,  $O((\log n)^{\log n}) = O(e^{\log n \log \log n})$ . Thus, we have  $O((\log n)!) = O(e^{\log n (\log \log n - 1)})$ .

Finally, the proxy and data aggregator again run the **Blind Ciphertext** and **Ciphertext Aggregation** phases to decrypt ciphertexts and do not make any modifications.

## 4 Protocol Analysis

### 4.1 Correctness

**Theorem 1.** *In our protocol, the data aggregator can accurately obtain the over-threshold set-union.*

*Proof.* W.l.o.g., we mainly take the  $t$  ciphertexts  $E_1, \dots, E_t$  as an example. We require to consider the following two cases.

**Case 1:** Assume the  $t$  ciphertexts encrypt the same data item  $x$ . Let  $\mathbf{v} = (x_{m-1}, \dots, x_0)$  denote the binary representation of  $x$ . For any  $i \in [t]$ , let  $\mathbf{x}_i = (r_{i,1}, h_{i,1}, 0^2, r_{i,2}, r_{i,3}, 0, 0, 0)_{\mathbb{B}_1}$  and  $\mathbf{y}_i = (r_{i,4}, r_{i,5}, 0^2, 0, 0, r_{i,6}, r_{i,7}, 0)_{\mathbb{B}_1}$ , we have  $F_{i,1} = \prod_{j=1}^t e(C_{i,1}, C_{j,2}) = \prod_{j=1}^t e(G_1, G_2)^{\mathbf{x}_i \cdot \mathbf{y}_j} = e(G_1, G_2)^{\sum_{j=1}^t (r_{i,1} r_{j,4} + h_{i,1} r_{j,5})}$ .

Similar to the above analysis, we have  $F_{i,2} = e(G_1, G_2)^{\sum_{j=1}^t (r_{i,1} r_{j,4} + \bar{h}_{j,1} \hat{r}_{i,3})}$ ,  $\hat{F}_{i,1} = e(G_1, G_2)^{\sum_{j=1}^t (r_{i,4} r_{j,1} + \bar{h}_{i,1} \hat{r}_{j,3})}$ ,  $\hat{F}_{i,2} = e(G_1, G_2)^{\sum_{j=1}^t (r_{i,4} r_{j,1} + h_{j,1} r_{i,5})}$ ,  $\bar{F}_i = e(G_1, G_2)^{\sum_{j=1}^t (h_{j,1} r_{i,5} + \bar{h}_{j,1} \hat{r}_{i,3}) + \sum_{j=1}^t \bar{r}_{j,4}}$ , where  $\mathbf{X}_{i,1} \cdot \mathbf{X}_{i,2} = \mathbf{v} \cdot \mathbf{1} + \mathbf{1} \cdot \mathbf{v} - 2\mathbf{v} \cdot \mathbf{v} = \sum_{l=0}^{m-1} x_l \oplus x_l = 0$ .  $H = \prod_{j=1}^t C_i = e(G_1, G_2)^{\sum_{j=1}^t (\beta_1 r_{j,5} + \gamma_1 \hat{r}_{j,3} + \bar{r}_{j,4})}$ .

According to the equations, we can deduce

$$\begin{aligned} \prod_{i=1}^t \left( \frac{F_{i,1} \hat{F}_{i,1} \bar{F}_i}{F_{i,2} \hat{F}_{i,2}} \right)^{\Delta_{\alpha_{i,1},B}(0)} &= \prod_{i=1}^t e(G_1, G_2)^{\Delta_{\alpha_{i,1},B}(0) \sum_{j=1}^t (f_{i,1} r_{j,5} + h_{i,1} \hat{r}_{j,3} + \bar{r}_{j,4})} \\ &= \prod_{i=1}^t e(G_1, G_2)^{\sum_{j=1}^t \Delta_{\alpha_{i,1},B}(0) (f_{i,1} r_{j,5} + h_{i,1} \hat{r}_{j,3} + \bar{r}_{j,4})} \\ &= e(G_1, G_2)^{\sum_{j=1}^t (\beta_1 r_{j,5} + \gamma_1 \hat{r}_{j,3} + \bar{r}_{j,4})} = H \end{aligned}$$

Then, we continue to compute  $\hat{U}$ . By using similar analysis, we can obtain

$$\begin{aligned} \hat{U} &= \left( \prod_{i=1}^t K_i \left( \frac{J_{i,2} \hat{J}_{i,2}}{J_{i,1} \hat{J}_{i,1} \bar{J}_i} \right)^{\Delta_{\alpha_{i,2},B}(0)} \right)^{\frac{1}{t}} = \left( e(G_1, G_2)^{tx} \right)^{\frac{1}{t}} = e(G_1, G_2)^x \\ U &= \log \hat{U} = x \end{aligned}$$

Therefore, we obtain the data item  $x$  which belongs to over-threshold set-union.

**Case 2:** Assume the  $t$  ciphertexts encrypt the different data items  $\{x_1, \dots, x_t\}$ . Let  $\mathbf{v}_i = (x_{i,l-1}, \dots, x_{i,0})$  denote the binary representation of  $x_i$ . Similar to the analysis in **Case 1**, we have

$$\bar{F}_i = \prod_{j=1}^t e(C_{i,5}, C_{j,6}) = e(G_1, G_2)^{\sum_{j=1}^t (h_{j,1} r_{i,5} + \bar{h}_{j,1} \hat{r}_{i,3} + \bar{r}_{j,4} + \zeta_{i,j})},$$

where  $\zeta_{i,j} = \bar{r}_{i,1} f_{a_1}(r) f_{b_1}(r) (\sum_{l'=0}^{l-1} x_{i,l'} \oplus x_{j,l'})$ . Thus, we can obtain

$$\prod_{i=1}^t \left( \frac{F_{i,1} \hat{F}_{i,1} \bar{F}_i}{F_{i,2} \hat{F}_{i,2}} \right)^{\Delta_{\alpha_{i,1},B}(0)} / H = e(G_1, G_2)^{\sum_{i=1}^t (\beta_1 r_{i,5} + \gamma_1 \hat{r}_{i,3} + \bar{r}_{i,4} + \sum_{j \in [t]} \zeta_{i,j})} / H = e(G_1, G_2)^{\sum_{i,j \in [t]} \zeta_{i,j}}$$

Since

$$0 \leq \sum_{i,j \in [t]} \zeta_{i,j} / f_{a_1}(r) f_{b_1}(r) = \sum_{i \in [t]} (\bar{r}_{i,1} (\sum_{l'=0}^{l-1} x_{i,l'} \oplus x_{j,l'})) \leq \max\{\bar{r}_{i,1} | i \in [t]\} tl < q,$$

the probability that  $\sum_{i,j \in [t]} \zeta_{i,j} = 0$  or  $e(G_1, G_2)^{\sum_{i,j \in [t]} \zeta_{i,j}} = 1$  is at most  $1/p^t$ , where  $f_{a_1}(r) f_{b_1}(r) \neq 0$ .

Therefore, by a union bound, our protocol makes wrong judgments or produces any false items with error probability at most  $\delta$ , where  $\delta \leq \frac{\binom{n}{t}}{p^t} \leq \frac{n^t}{t! p^t} \leq \frac{n^t}{n^{2t}} \leq \frac{1}{n^t}$ , ( $p > n^2$ ).

In summary, our protocol can accurately obtain the over-threshold set-union with high probability.

## 4.2 Privacy

We discuss the privacy of our protocol in this section. Since the private data items of participants are encrypted by using the distributed ElGamal encryption scheme, other participants and the proxy cannot decrypt them and obtain

the data privacy information based on the semantic security assumption. We only consider whether the aggregator can obtain extra information than what can be deduced from the over-threshold set-union or not. According to Theorem 1, we know that if  $t$  ciphertexts are generated by using the same data item, we can obtain the data item. Thus, we only consider the ciphertexts whose corresponding plaintexts do not contain  $t$  or more the same data item. We first consider a special case: given any  $x, y_0, y_1$  (corresponding vectors  $\mathbf{x}, \mathbf{y}_0, \mathbf{y}_1$ ),  $\{\overline{C}_i\}_{i=1}^t = \{C_i, \{C_{i,j}\}_{j=1}^6\}_{i=1}^t$  and  $\{\overline{D}_i\}_{i=1}^t = \{D_i, \{D_{i,j}\}_{j=1}^6\}_{i=1}^t$ , where  $\{\overline{C}_2, \dots, \overline{C}_t\}, \{\overline{D}_2, \dots, \overline{D}_t\}$  are the ciphertexts of  $\mathbf{y}$ , and  $\overline{C}_1, \overline{D}_1$  are the ciphertexts of  $\mathbf{y}_0$  or  $\mathbf{y}_1$ .

Without loss of generality, we use  $\mathbf{Y}_{i,1}$  and  $\mathbf{Y}_{i,2}$  to denote the vectors  $(\mathbf{y}_i, \mathbf{1}, \mathbf{y}_i)$  and  $(\mathbf{1}, \mathbf{y}_i, -2\mathbf{y}_i)$ , respectively. Then, we define a series of games as follows.

**Game  $\mathbf{G}_0$**  : This game is a original on where  $\overline{D}_1$  is the ciphertext of  $y_0$ .

**Game  $\mathbf{G}_1$**  : This game is the same as  $\mathbf{G}_1$  except that  $D_{1,1} = (s'_{1,1} \cdot \mathbf{Y}_{1,1}, s'_{1,2}, h_{1,2}, 0^{3l+2}, s_{1,3}, s_{1,4}, 0, 0, 0)_{\mathbb{C}_1}$ , which is denoted by  $D'_{1,1}$ .

**Game  $\mathbf{G}_2$**  : This game is the same as  $\mathbf{G}_1$  except that  $D_{1,2} = (s'_{1,5} \cdot \mathbf{Y}_{1,2}, s'_{1,6}, s'_{1,7}, 0^{3l+2}, s_{1,8}, s_{1,9}, 0, 0, 0)_{\mathbb{C}_1^*}$ , which is denoted by  $D'_{1,2}$ .

**Game  $\mathbf{G}_3$**  : This game is the same as  $\mathbf{G}_2$  except that  $D_{1,3} = (s'_{1,1} \cdot \mathbf{Y}_{1,1}, s'_{1,6}, \overline{h}_{1,2}, 0^{3l+2}, \hat{s}_{1,2}, \hat{s}_{1,3}, 0, 0, 0)_{\mathbb{C}_1}$ , which is denoted by  $D'_{1,3}$ .

**Game  $\mathbf{G}_4$**  : This game is the same as  $\mathbf{G}_3$  except that  $D_{1,4} = (s'_{1,4} \cdot \mathbf{Y}_{1,2}, s'_{1,2}, s'_{1,5}, 0^{3l+2}, \hat{s}_{1,6}, \hat{s}_{1,7}, 0, 0, 0)_{\mathbb{C}_1^*}$ , which is denoted by  $D'_{1,4}$ .

**Game  $\mathbf{G}_5$**  : This game is the same as  $\mathbf{G}_4$  except that  $D_{1,5} = (\overline{s}'_{1,1} \cdot \mathbf{Y}_{1,1}, s'_{1,7}, \hat{s}'_{1,5}, 1, 0^{3l+3}, \overline{s}_{1,3}, \overline{s}_{1,4}, 0, 0, 0)_{\mathbb{C}_1}$ , which is denoted by  $D'_{1,5}$ .

**Game  $\mathbf{G}_6$**  : This game is the same as  $\mathbf{G}_5$  except that  $D_{1,6} = (\overline{s}'_{1,4} \cdot \mathbf{Y}_{1,2}, h_{1,2}, \overline{h}_{1,2}, \overline{s}'_{1,5}, 0^{3l+3}, \overline{s}_{1,8}, \overline{s}_{1,9}, 0, 0, 0)_{\mathbb{C}_1^*}$ , which is denoted by  $D'_{1,6}$ .

**Game  $\mathbf{G}_7$**  : This game is the same as  $\mathbf{G}_6$  except that  $C_1 = e(G_1, G_2)^{y_1 + \beta_2 r'_{i,7} + \gamma_2 \overline{r}'_{i,5} + \overline{r}'_{i,5}}$  which is denoted by  $C'_1$ .

**Lemma 1** *Given  $\{(q, \mathbb{G}_1, \mathbb{G}_2, \mathbb{G}_T, e), \widehat{\mathbb{B}}, \widehat{\mathbb{B}}^*\}$ , when the SXDH and XDLIN assumptions hold, for any PPT adversary,  $D_{1,1} \approx_c D'_{1,1}$ , where  $\widehat{\mathbb{B}} = (\mathbf{b}_{2,1}, \dots, \mathbf{b}_{2,6l+4}, \mathbf{b}_{2,6l+5}, \mathbf{b}_{2,6l+6}), \widehat{\mathbb{B}}^* = (\mathbf{b}_{2,1}^*, \dots, \mathbf{b}_{2,6l+4}^*, \mathbf{b}_{2,6l+7}^*, \mathbf{b}_{2,6l+8}^*)$  and  $\approx_c$  denotes the computationally indistinguishable.*

The lemma can be proved by using similar methods in [3] and [16], which we omit here.

**Lemma 2** *Suppose that the XDLIN assumption holds in  $\mathbb{G}_1$  and  $\mathbb{G}_2$ , for any PPT adversary, **Game  $\mathbf{G}_0$**   $\approx_c$  **Game  $\mathbf{G}_1$***

*Proof.* Invoking Lemma 1, we have  $D_{1,1} \approx_c D'_{1,1}$  and  $(D_{1,1}, D_{j,1}) \approx_c (D'_{1,1}, D_{j,1})$  for any  $j \in \{1, 2, \dots, t\}$ . Let  $\zeta_0 = x \oplus y_0$ ,  $\zeta'_0 = x \oplus y_1$  and  $\zeta''_0 = y_0 \oplus y_1$ . Then

$$\begin{aligned} J_{1,1} &= \prod_{j=1}^t e(D_{1,1}, D_{j,2}) = e(G_1, G_2)^{\sum_{j=1}^t (s_{1,2}s_{j,6} + h_{1,2}s_{j,7}) + \mu} \\ J'_{1,1} &= \prod_{j=1}^t e(D'_{1,1}, D_{j,2}) = e(G_1, G_2)^{\sum_{j=1}^t (s'_{1,2}s_{j,6} + h_{1,2}s_{j,7}) + \mu'}, \\ J_{i,1} &= \prod_{j=1}^t e(D_{i,1}, D_{j,2}) = e(G_1, G_2)^{\sum_{j=1}^t (s_{i,2}s_{j,6} + h_{i,2}s_{j,7}) + s_{i,1}s_{j,5}\zeta_0} (i \neq 1), \\ J'_{i,1} &= \prod_{j=1}^t e(D_{i,1}, D_{j,2}) = J_{i,1} (i \neq 1). \\ \hat{J}'_{1,2} &= e(D_{1,2}, D'_{1,1}) \prod_{j=2}^t e(D_{1,2}, D_{j,1}) = e(G_1, G_2)^{\sum_{j=2}^t (s_{1,6}s_{j,2} + s_{1,7}h_{j,2}) + \mu''}, \\ \hat{J}'_{i,2} &= e(D_{i,2}, D'_{1,1}) \prod_{j=2}^t e(D_{i,2}, D_{j,1}) = e(G_1, G_2)^{\sum_{j=2}^t (s_{1,6}s'_{j,2} + s_{1,7}h_{j,2}) + \mu'''} (i \neq 1), \end{aligned}$$

where  $\mu = \sum_{i=2}^t s_{1,1}s_{i,5}\zeta_0$ ,  $\mu' = s'_{1,1}s_{1,5}\zeta''_0 + \sum_{i=2}^t s'_{1,1}s_{i,5}\zeta'_0$ ,  $\mu'' = s_{1,6}s'_{1,2} + h_{1,2}s_{1,7} + s'_{1,1}s_{1,5}\zeta''_0 + \sum_{i=2}^t s_{i,1}s_{1,5}\zeta_0$  and  $\mu''' = s_{i,6}s'_{1,2} + h_{1,2}s_{i,7} + s_{i,5}s'_{1,1}\zeta'_0$ .

$$M' = \prod_{i=1}^t \left( \left( \frac{J_{i,1}\hat{J}_{i,1}\bar{J}_i}{J_{i,2}\hat{J}_{i,2}} \right)^{\Delta_{\alpha_{i,2},B}(0)} / H_i \right) = e(G_1, G_2)^{a_1 - y_0 - (t-1)x}$$

where  $a_1 = \Delta_{\alpha_{1,2},B}(0) \left( \zeta_0 \sum_{j=2}^t (s_{1,1}s_{j,5} - \hat{s}_{1,4}\hat{s}_{j,1} + \hat{s}_{1,1}\hat{s}_{j,4} - s_{1,5}s_{j,1} + \bar{s}_{1,1}\bar{s}_{j,4}) \right) + \zeta_0 \sum_{i=2}^t \Delta_{\alpha_{i,2},B}(0) \left( s_{i,1}s_{1,5} - \hat{s}_{i,4}\hat{s}_{1,1} + \hat{s}_{i,1}\hat{s}_{1,4} - s_{i,5}s_{1,1} + \bar{s}_{i,1}\bar{s}_{1,4} \right)$ . Similarly, we have

$$M'' = \prod_{i=1}^t \left( \left( \frac{J'_{i,1}\hat{J}_{i,1}\bar{J}_i}{J_{i,2}\hat{J}'_{i,2}} \right)^{\Delta_{\alpha_{i,2},B}(0)} / H_i \right) = e(G_1, G_2)^{a'_1 - y_0 - (t-1)x}$$

where  $a'_1 = \sum_{i=2}^t \Delta_{\alpha_{i,2},B}(0) \left( \zeta_0 (s_{i,1}s_{1,5} - \hat{s}_{i,4}\hat{s}_{1,1} + \bar{s}_{i,1}\bar{s}_{1,4} + \hat{s}_{i,1}\hat{s}_{1,4}) - \zeta'_0 s_{i,5}s'_{1,1} \right) + \Delta_{\alpha_{1,2},B}(0) \left( \sum_{i=2}^t \zeta'_0 s'_{1,1}s_{i,5} + \zeta''_0 s'_{1,1}s_{1,5} - \sum_{i=2}^t \zeta_0 \hat{s}_{1,4}\hat{s}_{i,1} + \sum_{i=2}^t \zeta_0 \hat{s}_{i,4}\hat{s}_{1,1} - \sum_{i=2}^t \zeta_0 s_{1,5}s_{i,1} - \zeta''_0 s'_{1,1}s_{1,5} + \sum_{i=2}^t \zeta_0 \bar{r}_{1,1}\bar{r}_{i,4} \right) + \Delta_{\alpha_{1,2},B}(0) \sum_{j \in [t]} (s'_{1,2}s_{j,6} - s_{1,2}s_{j,6} + s_{1,6}s_{j,2} - s_{1,6}s'_{j,2}) + \sum_{i=2}^t \Delta_{\alpha_{i,2},B}(0) (s_{i,6}s_{1,2} - s_{i,6}s'_{1,2})$ .

Since  $s_{1,1}, s'_{1,1} \xleftarrow{U} Z_p$ , we have  $M' \approx_c M''$ . Thus, we obtain  $\{D_1, D_{1,1}, D_{1,2}, \dots, D_{1,6}\} \approx_c \{D_1, D'_{1,1}, D_{1,2}, \dots, D_{1,6}\}$ , i.e., **Game**  $\mathbf{G}_0 \approx_c \mathbf{Game} \mathbf{G}_1$ .

By using a similar approach as in Lemma 2, we have **Game**  $\mathbf{G}_i \approx_c \mathbf{Game} \mathbf{G}_{i+1}$  for any  $i \in \{1, \dots, 6\}$ .

It is noted that the ciphertexts in **Game G<sub>7</sub>** are about  $y_1$ . Thus, Lemma 3 can be proved by using the above lemmas.

**Lemma 3** *Suppose that the SXDH and XDLIN assumptions hold in  $\mathbb{G}_1$  and  $\mathbb{G}_2$ , given any  $x, y_0, y_1$  (corresponding vectors  $\mathbf{x}, \mathbf{y}_0, \mathbf{y}_1$ ) and  $\{\bar{D}_i\}_{i=1}^t = \{D_i, \{D_{i,j}\}_{j=1}^6\}_{i=1}^t$ , the advantage of any PPT adversary which guesses to  $b \in \{0, 1\}$  is negligible, where  $\{\bar{D}_2, \dots, \bar{D}_t\}$  are the ciphertexts of  $\mathbf{y}$  and  $\bar{D}_1$  is the ciphertexts of  $\mathbf{y}_0$  or  $\mathbf{y}_1$ .*

For the ciphertexts  $\{\bar{C}_i\}_{i=1}^t = \{C_i, \{C_{i,j}\}_{j=1}^6\}_{i=1}^t$  of  $y_0$  and  $x$ , similar to the above analysis, we have

$$N' = \prod_{i=1}^t \left( \frac{F_{i,1} \hat{F}_{i,1} \bar{F}_i}{F_{i,2} \hat{F}_{i,2}} \right)^{\Delta_{\alpha_{i,1}, B}(0)} = e(G_1, G_2)^{a'_1}$$

where  $a'_1 = f_{a_1}(r) f_{b_1}(r) \left( \Delta_{\alpha_{1,1}, B}(0) \bar{r}_{1,1} (t-1) \zeta_0 + \sum_{i=2}^t \Delta_{\alpha_{i,1}, B}(0) \bar{r}_{i,1} \zeta_0 \right)$ . Although  $\bar{r}_{i,1} (\in Z_{q'})$  does not be chosen uniformly at random from  $Z_p$  for any  $i \in [t]$ ,  $\zeta_0$  is still secure if  $\bar{r}_{i,1}$  sufficiently bigger than  $\zeta_0$  [?]. Thus, we have the following lemma.

**Lemma 4** *Suppose that the SXDH and XDLIN assumptions hold in  $\mathbb{G}_1$  and  $\mathbb{G}_2$ , given any  $x, y_0, y_1$  (corresponding vectors  $\mathbf{x}, \mathbf{y}_0, \mathbf{y}_1$ ), the ciphertexts  $\{\bar{C}_i\}_{i=1}^t = \{C_i, \{C_{i,j}\}_{j=1}^6\}_{i=1}^t$  do not reveal information about the  $x, y_0$  or  $x, y_1$ , where  $\{\bar{C}_2, \dots, \bar{C}_t\}$  are the ciphertexts of  $\mathbf{y}$  and  $\bar{C}_1$  is the ciphertexts of  $\mathbf{y}_0$  or  $\mathbf{y}_1$ .*

**Lemma 5** *Suppose that the XDLIN assumption holds in  $\mathbb{G}_1$  and  $\mathbb{G}_2$ , given any  $x_{0,1}, \dots, x_{0,t}, x_{1,1}, \dots, x_{1,t}$  (corresponding vectors  $\mathbf{x}_{0,1}, \dots, \mathbf{x}_{0,t}, \mathbf{x}_{1,1}, \dots, \mathbf{x}_{1,t}$ ) and ciphertexts  $\{\bar{D}_i\}_{i=1}^t$ , the advantage of any PPT adversary which guesses to  $b \in \{0, 1\}$  is negligible, where  $\{\bar{D}_1, \dots, \bar{D}_t\}$  are the ciphertexts of  $\{\mathbf{x}_{0,1}, \dots, \mathbf{x}_{0,t}\}$  or  $\{\mathbf{x}_{1,1}, \dots, \mathbf{x}_{1,t}\}$ , where  $\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,t}$  are not the same vector for  $i \in \{0, 1\}$ .*

*Proof.* We define a series of games and use a hybrid argument to prove this lemma as follows.

**Game H<sub>0</sub>** : The ciphertexts are  $\{\bar{D}_i(x_{0,i})\}_{i=1}^t$ .

**Game H<sub>1</sub>** : This game is the same as **Game H<sub>0</sub>** except that the first ciphertexts are  $\bar{D}_i(x_{1,1})$ , i.e.,  $\{\bar{D}_i(x_{1,1}), \bar{D}_i(x_{0,2}), \dots, \bar{D}_i(x_{0,t})\}$

...

**Game H<sub>i</sub>** : This game is the same as **Game H<sub>i-1</sub>** except that the  $i$ -th ciphertexts are  $\bar{D}_i(x_{1,i})$ , i.e.,  $\{\bar{D}_i(x_{1,1}), \dots, \bar{D}_i(x_{1,i}), \dots, \bar{D}_i(x_{0,t})\}$

...

**Game H<sub>t</sub>** : This game is the same as **Game H<sub>t-1</sub>** except that the  $i$ -th ciphertexts are  $\bar{D}_i(x_{1,t})$ , i.e.,  $\{\bar{D}_i(x_{1,1}), \dots, \bar{D}_i(x_{1,i}), \dots, \bar{D}_i(x_{1,t})\}$

By using Lemma 3, we have **Game H<sub>i</sub>**  $\approx_c$  **Game H<sub>i+1</sub>** for any  $i \in \{0, \dots, t-1\}$ . So we have **Game H<sub>0</sub>**  $\approx_c$  **Game H<sub>t</sub>**, i.e.,  $\{\bar{D}_i(x_{0,i})\}_{i=1}^t \approx_c \{\bar{D}_i(x_{1,i})\}_{i=1}^t$ .

Similar to the analysis in Lemmas 4 and 5, we have the following lemma.

**Lemma 6** *Suppose that the SXDH and XDLIN assumptions hold in  $\mathbb{G}_1$  and  $\mathbb{G}_2$ , given any  $x_{0,1}, \dots, x_{0,t}, x_{1,1}, \dots, x_{1,t}$  (corresponding vectors  $\mathbf{x}_{0,1}, \dots, \mathbf{x}_{0,t}, \mathbf{x}_{1,1}, \dots, \mathbf{x}_{1,t}$ ), the ciphertexts  $\{\bar{C}_i\}_{i=1}^t = \{C_i, \{C_{i,j}\}_{j=1}^6\}_{i=1}^t$  do not reveal information about the  $x_{0,j}, x_{0,i}$  or  $x_{1,j}, x_{1,i}$  for any  $i, j \in \{1, \dots, t\}$  and  $i \neq j$ , where  $\{\bar{C}_1, \dots, \bar{C}_t\}$  are the ciphertexts of  $\{\mathbf{x}_{0,1}, \dots, \mathbf{x}_{0,t}\}$  or  $\{\mathbf{x}_{1,1}, \dots, \mathbf{x}_{1,t}\}$ , where  $\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,t}$  are not the same vector for  $i \in \{0, 1\}$ .*

**Theorem 2.** *Our protocol is secure under the SXDH and XDLIN assumptions.*

*Proof.* According to lemmas 5 and 6, we have the ciphertexts  $\{\bar{C}_i\}_{i \in \mathcal{A}}$  and  $\{\bar{D}_i\}_{i \in \mathcal{A}}$  do not reveal privacy information about the participants' data items. Thus, our protocol is secure.

### 4.3 Complexity

It is not difficult to see that the communication complexity of **System Setup**, **Ciphertext Generation** and **Blind Ciphertext** of our protocol are  $O(n(\log n)^2)$ ,  $O(n \log n \log p)$  and  $O(n \log n \log p)$ , respectively. Thus, the total communication complexity of our protocol is  $O(n(\log n)^2)$ , where  $p > n^2$ .

## 5 Conclusions

In this paper, we propose a novel communication-efficient and privacy-preserving protocol for the over-threshold set-union problem. Compared with previous schemes, our protocol has nearly optimal communication complexity and meanwhile can protect data privacy.

## References

1. Abe, M., Chase, M., David, B., Kohlweiss, M., Nishimaki, R., Ohkubo, M.: Constant-size structure-preserving signatures: Generic constructions and simple assumptions. *Journal of Cryptology* **29**(4), 833–878 (2016)
2. Applebaum, B., Ringberg, H., Freedman, M.J., Caesar, M., Rexford, J.: Collaborative, privacy-preserving data aggregation at scale. In: *PETS*. pp. 56–74. Springer (2010)
3. Bishop, A., Jain, A., Kowalczyk, L.: Function-hiding inner product encryption. In: *ASIACRYPT*. pp. 470–491. Springer (2015)
4. Bonawitz, K., Ivanov, V., Kreuter, B., Marcedone, A., McMahan, H.B., Patel, S., Ramage, D., Segal, A., Seth, K.: Practical secure aggregation for privacy-preserving machine learning. In: *ACM CCS*. pp. 1175–1191 (2017)
5. Burkhart, M., Dimitropoulos, X.: Fast privacy-preserving top-k queries using secret sharing. In: *ICCCN*. pp. 1–7. IEEE (2010)
6. Castelluccia, C., Chan, A.C., Mykletun, E., Tsudik, G.: Efficient and provably secure aggregation of encrypted data in wireless sensor networks. *ACM TOSN* **5**(3), 1–36 (2009)



7. Gennaro, R., Jarecki, S., Krawczyk, H., Rabin, T.: Secure distributed key generation for discrete-log based cryptosystems. *Journal of Cryptology* **20**(1), 51–83 (2007)
8. He, W., Liu, X., Nguyen, H., Nahrstedt, K., Abdelzaher, T.: Pda: Privacy-preserving data aggregation in wireless sensor networks. In: *INFOCOM*. pp. 2045–2053. IEEE (2007)
9. Kim, M., Mohaisen, A., Cheon, J.H., Kim, Y.: Private over-threshold aggregation protocols over distributed datasets. *IEEE Transactions on Knowledge and Data Engineering* **28**(9), 2467–2479 (2016)
10. Kissner, L., Song, D.: Privacy-preserving set operations. In: *CRYPTO*. pp. 241–257. Springer (2005)
11. Li, Q., Cao, G.: Efficient and privacy-preserving data aggregation in mobile sensing. In: *ICNP*. pp. 1–10. IEEE (2012)
12. Luby, M., Rackoff, C.: How to construct pseudorandom permutations from pseudorandom functions. *SIAM Journal on Computing* **17**(2), 373–386 (1988)
13. Menezes, A.J., Katz, J., Van Oorschot, P.C., Vanstone, S.A.: *Handbook of applied cryptography*. CRC press (1996)
14. Okamoto, T., Takashima, K.: Fully secure functional encryption with general relations from the decisional linear assumption. In: *CRYPTO*. pp. 191–208. Springer (2010)
15. Shamir, A.: How to share a secret. *Communications of the ACM* **22**(11), 612–613 (1979)
16. Tomida, J., Abe, M., Okamoto, T.: Efficient functional encryption for inner-product values with full-hiding security. In: *ISC*. pp. 408–425. Springer (2016)
17. Woodruff, D.P., Zhang, Q.: When distributed computation is communication expensive. *Distributed Computing* **30**(5), 309–323 (2017)