

# Private Top- $\kappa$ Aggregation Protocols<sup>\*</sup>

Myungsun Kim<sup>1</sup>, Abedelaziz Mohaisen<sup>2</sup>, Jung Hee Cheon<sup>3</sup>, and Yongdae Kim<sup>4</sup>

<sup>1</sup> Information Security Engineering Department  
University of Suwon  
`msunkim@suwon.ac.kr`

<sup>2</sup> Verisign Labs, VA, USA  
`amohaisen@verisign.com`

<sup>3</sup> Department of Mathematical Sciences  
Seoul National University  
`jhcheon@snu.ac.kr`

<sup>4</sup> Department of Electrical Engineering  
Korea Advanced Institute of Science and Technology  
`yongdaek@ee.kaist.ac.kr`

**Abstract.** In this paper, we revisit the private top- $\kappa$  data aggregation problem. First we formally define the problem’s security requirements as both data and user privacy goals. To achieve both goals, and to strike a balance between efficiency and functionality, we devise a novel cryptographic construction that comes in two schemes; a fully decentralized simple construction and its practical and semi-decentralized variant. Both schemes are provably secure in the semi-honest model. We analyze the computational and communication complexities of our construction, and show that it is much more efficient than the existing protocols in the literature.

**Keywords:** Privacy-preservation, top- $\kappa$  items, data privacy, user privacy

## 1 Introduction

Privacy has become a big concern in many applications of computing in various domains, and for various technical, technological, and ethical reasons. For example, privacy in collaborative data aggregation—an application that arises in sensor data [18] as well as Internet data aggregation and analysis [4]—has arisen for both maintaining users privacy and business secrecy; both ethical and technological reasons. Privacy when dealing with medical records arises due to both legislations and ethical concerns [30]. Privacy in e-voting [21,10,24] arises due to the ethical essence of voting systems; users cast votes and their identities are to be kept private and unlinkable to their votes. In the vast majority of these applications, and despite great differences in the domains where such applications are applied, many commonalities exist among them in the way they operate. Such commonalities are seen in the underlying computing primitives utilized over the data generated by such applications. Thus, providing an overall privacy-preserving application would suffice by augmenting such primitives to be privacy-preserving as building blocks.

Of particular interest in many applications is the problem of *computing the top- $\kappa$  elements* in a *private* manner. While most of the aforementioned applications may require such primitive, a typical scenario of an application that is easy to explain, which involves such primitive, is network traffic distribution. In a typical scenario of this application,  $n$  network sensors need to jointly analyze the security alert broadcasted by different sources in order to find the top- $\kappa$  suspect sites. In such an application, and without losing generality, each of such sensors—supposedly under control of different authorities—has a set of suspects and would like to collaboratively compute the top most frequent on each of these sets (say  $\kappa$  of them) without revealing the set of suspects to other sensors with whom she collaborates.

In this paper, we consider the problem of private top- $\kappa$  computing, and introduce a construction and its formal security analysis. In the main protocol, our operation setting is *fully*

---

<sup>\*</sup> This work was partly done while M. Kim and J. Cheon were visiting the University of Minnesota, and while A. Mohaisen and Y. Kim were with the University of Minnesota.

*decentralized* among  $n$  users over non-partitioned data. For efficiency, we refrain from using secure multiparty computation, and construct an efficient private top- $\kappa$  protocol which is both data-private and user-private. Our construction strikes a balance between the consumed resources and achieved security. It also satisfies both privacy requirements formalized as goals, thus outperforming existing literature. In the following, we start by formally defining the problem, then summarizing our results, and finally outlining the contribution and organization of this work.

## 1.1 Problem definition

Let there be  $n$  users denoted by  $u_i, 1 \leq i \leq n$ , and each of them has a private (multi-)set  $X_i$  of cardinality  $k$ . For simplicity, assume that the cardinality of each multiset is equal to each other. (We can efficiently handle the case where the cardinality of all multisets are different from each other by adding random elements.) There may exist one or more elements such that  $\alpha_{i,j} = \alpha_{i,j'}$  for some  $j \neq j'$ .

**PRIVATE TOP- $\kappa$  AGGREGATION PROBLEM:** By the multiplicity of an element of a multiset we mean the number of times it appears in the multiset. Let  $\kappa \in \mathbb{N}$ , and assume  $\kappa$  has been implicitly predefined among all users. Then the problem at hand is defined as follows: Given  $n$  multisets of cardinality  $k$ , find a set  $Z = \{\alpha_1, \dots, \alpha_\kappa\} \subset \mathcal{U} = \bigcup_{i=1}^n X_i$  such that (i) for all elements  $\alpha \in \mathcal{U}$ , if  $\alpha$  has a multiplicity greater than or equal to  $\kappa$ , then  $\alpha \in Z$ , (ii) no polynomial-time algorithm can learn any element other than the output of a top- $\kappa$  protocol, and (iii) no polynomial-time algorithm should know which output of the execution belongs to which user.

With such definition of the problem in mind, we now proceed to outline other possible approaches taken to solve this problem. One trivial and straightforward technique to solve the problem is to use a trusted third party (TTP). In such case, each user sends his private set to such TTP, which then performs the top- $\kappa$  aggregation task and reports the result back to each user. However, it is well known that finding such TTP is not always possible in many real-world applications. Moreover, even though it is feasible to provide TTP, compromising the TTP could lead to a complete privacy loss for all participating users.

We can also base a top- $\kappa$  protocol on secure multiparty computation (SMC). SMC allows us to securely compute a function over private data such that users only learn the result of the function and nothing else. However, despite recent advances in their efficiency, SMC methods still require substantial computation and communication costs which make them impractical for real-world applications than mainly operate on large datasets.

A final approach is realizable using existing private set-operation protocols such as [22,28,19], especially multiset union protocols. These protocols securely compute all elements appearing in the union of input multisets at least  $\tau$  times. Here all of them demand a priori-knowledge of the threshold value  $\tau$ . However, such parameter depends on the data itself, and might not be available in advance. Furthermore, computing such parameter across different users while maintaining the privacy of the data at the side of every user is obviously a non-trivial task.

**Remark 1** *There have been many results in the literature with titles containing the term “top- $\kappa$ ”, which could confuse the reader as to the extent of the novelty of our work—examples of such literature include [29,31,33,34]. We stress that these protocols produce the greatest  $\kappa$  elements in the union of the given sets (or multisets) in a private manner, and thus are different in their end results from our protocol. In particular, these schemes take as an input numerical-valued sets. For example, there is a secure method for finding the  $\kappa$ -th ranked element in multiple multisets by Aggarwal et al. [2]. Repeatedly applying this protocol we can efficiently find the biggest  $\kappa$  elements in a distributed list. Hence our protocol is essentially different from these protocols.*

## 1.2 Our Approach—Informal Descriptions

We observe that the problem of finding the most frequent  $\kappa$  elements from the given  $n$  multisets is conceptually similar to finding most popular  $\kappa$  candidates from  $n$  candidates (as it is the typical scenario in a voting system). Motivated by that, we borrow ideas from voting protocols

to solve this problem—see [21,10,12,24,16] and refer to [1] for a detailed exposition on the voting problem and related solutions.

The most non-trivial part of top- $\kappa$  protocols is that we should satisfy two privacy requirements, namely the data privacy and user-privacy, at the same time. First let us take a closer look at e-voting protocols. In e-voting protocols, each ballot is mixed with a shuffle scheme which plays a crucial role in removing linkability between voters and ballots, which would hint user privacy. In fact, in e-voting literature this notion is called unlinkability. In order to emphasize the difference with e-voting protocols, we use the term user privacy. Now assuming that each element in multisets is encrypted and shuffled as in e-voting protocols, all encrypted elements can be decrypted, especially in a threshold manner, while satisfying user privacy. However, all non-top- $\kappa$  elements also are revealed, which violates data privacy in our application. Thus we need a way to keep data privacy even when all encrypted elements were decrypted.

For this purpose we introduce an efficient function  $E$  that commutes with an underlying public-key encryption. More specifically, let  $\text{Enc}$  be a public-key encryption algorithm and  $\text{Dec}$  be the corresponding decryption algorithm. We demand that

- (i) for all  $s$  and for all  $pk$ ,  $E_s \circ \text{Enc}_{pk} = \text{Enc}_{pk} \circ E_s$
- (ii) for all elements  $\alpha$ , given  $\text{Dec}_{sk}(\text{Enc}_{pk}(E_s(\alpha)))$  no algorithm can efficiently find  $\alpha$  without a secret information  $s$ .

We call this notion *double encryption*.<sup>5</sup> In conclusion, our main technique is to shuffle doubly encrypted elements by each user. We should notice that all shuffle algorithms used in e-voting protocols rely on the re-randomizable property of underlying homomorphic encryption (e.g. see [12,24,25,17,16]), but its re-randomization algorithm does not change the plaintexts of input ciphertexts. However, in our protocol a double encryption scheme will change the plaintexts of input ciphertexts, which is the main difference from existing shuffle algorithms. To this end, in the following we summarize our results and findings in this paper.

### 1.3 Summary of Our Results

In this paper, we present a formal definition of private top- $\kappa$  protocol and its security. Our operation setting is *fully decentralized* among  $n$  users over non-partitioned data. For the efficiency of our protocol, we refrain from using secure multiparty computation and construct an efficient private top- $\kappa$  protocol which is both data-private and user-private. Our construction strikes a real balance in the consumed resources and achieved security, and satisfies both privacy requirements. In particular, in its efficiency, our construction is comparable to the work in [4], which achieves its efficiency by giving up decentralized communication model, and in its security guarantees it is comparable to the work [6], which is secure but resources exhaustive. Our protocol on the other hand overcome the limitations and shortcomings of those protocols.

More specifically, our scheme does not requires a trusted party to set up the keys. Note that using Paillier cryptosystem [27] in a threshold manner demands some trusted setup. Moreover, our proposed protocol has several desirable features as follows:

- (1) It has  $\mathcal{O}(n^2k)$  computational complexity where  $n$  is the number of users and  $k$  is the cardinality of each user's set, assuming  $\kappa \leq k$ .
- (2) It has  $\mathcal{O}(n^2k)$  communication complexity.
- (3) It has a linear round complexity in the number of users.

In general, real-world applications has  $n$  much smaller than  $k$ , which further justifies the efficiency of our protocol. This is, our protocol is beneficial in such real-world applications where the number of participating users is small but the size of their multisets is large. It remains an important open problem to devise a protocol whose round complexity does not depend on the number of users. Then we could make our protocol have  $\mathcal{O}(nk)$  computation and communication complexity.

<sup>5</sup> Applebaum et al. [4] also use the same term but differently; in their scheme double encryption means that a message  $a$  is encrypted by two public keys.

**Organization.** The rest of this paper is organized as follows. In Section 2, we discuss the related work with extensive analysis and comparison to our work. In Section 3, we outline the preliminaries required for understanding the rest of the paper, including double encryption, our formalism of top- $\kappa$  protocol and its security, and cryptographic primitives used in the context of computing the top- $\kappa$ . In Section 4, we introduce our construction that comes into two forms with different requirements and guarantees and meets data privacy and user privacy. In Section 5 we analyze the security and complexity of our work, by proving its security and showing its resources consumption requirements. Remarks on extending our protocol to be secure in the malicious model are in Section 6. Finally, we draw concluding remarks and point future work in Section 7.

## 2 Related Work

There has been plenty of work in the literature to solve the problem of private data aggregation. Such schemes can be classified under three schools of thoughts: fully centralized, fully decentralized, and semi-decentralized. While the centralized schemes assume the existence of a trusted third party (TTP), which makes such schemes of less interest from the cryptographic and practical point of views—as such TTP is not immune to compromise and malicious behavior, the fully decentralized schemes utilize cryptographic primitives and protocols to replace the centralized TTP. On the other hand, semi-decentralized schemes try to bridge the functional and security gap between both of the aforementioned schools of thoughts. As they are of particular relevance to our work, in the following we elaborate on the two latter schools of thoughts: the decentralized and semi-decentralized private data aggregation protocols.

As we mentioned earlier, decentralized solutions to the problem try to replace the centralized TTP with cryptographic constructions, which comes in different forms leading to several directions of research. One direction is based on SMC, as it is the case in [6]. However, at the core of that protocol’s drawbacks is its inefficiency: since it uses Yao’s garbled circuits [32], the computational resources required for ordinary data sizes is overwhelmingly high. Furthermore, as the datasets become disjoint, the efficiency of such construction decreases sharply. In [6], Burkhart and Dimitropoulos—in what we consider to be the most relevant work to ours—have devised a construction in which *the round complexity is linear to the number of bits* in the data elements. However, due to using sketches to improve the efficiency of subprotocols, the aggregate results are *probabilistic*. Furthermore, while the work in [6] is efficient in terms of its computational complexity, this efficiency comes at the expense of high round complexity.

Although not addressing the same problem at hand, Kissner and Song [22] devised an over-threshold set union protocol—where *a threshold value  $\tau$  should be given in advance*—to find all elements appearing in the union of input multisets at least  $\tau$  times. The protocol requires a priori knowledge of the threshold, although operates in a decentralized manner. We compare it to our work (in Section 5) as a classical related work.

The final school of thoughts to solving the problem uses *semi-decentralized* constructions and is represented by the recent work of Applebaum et al. in [4]. In that work, the authors aim to enhance the efficiency of fully-decentralized instantiations by adding new entities: proxy and database (DB). However the proxy and the DB are assumed to be semi-honest restricting the possibility of coalition between proxy and DB. This allows to obtain a constant round protocol. While the authors claim that both proxy and DB are expected to act as semi-honest, we disagree with the assumption and believe that this assumption is too strong both from a theoretical and a practical point of view. Furthermore, their scheme extensively relies on oblivious transfer (OT) [23], which is a very expensive public-key primitive since it may require two modular exponentiations per invocation, and run for each bit of the user’s data element. For efficiency, their scheme may use an efficient version of a standard OT [20], but whose security is in random oracle model [5].

To sum up, Table 1 summarizes properties and efficiency of existing solutions compared with our proposed protocol. Computational complexity is expressed in terms of the number of multiplications over modulo  $p$ . For simplicity, we assume that multisets have values less than the prime  $p$ . Note that Applebaum et al.’s protocol requires both ElGamal encryption [11] and

Goldwasser-Micali encryption [14], but we assume that both encryption systems use the same size modulus.

**Table 1.** Summary and Comparison

	Comm. Model	Round Cpx	Comp. Cpx	Comm. Cpx
Ours	Fully decentralized	$\mathcal{O}(n)$	$\mathcal{O}(n^2 k \log p)$	$\mathcal{O}(n^2 k \log p)$
[6]	Fully decentralized	$\mathcal{O}(n(n + k \log k) \log p)$	$\mathcal{O}(n^2 k)$	$\mathcal{O}(n^2 k \log p)$
[4]	Semi-decentralized	$\mathcal{O}(1)$	$\mathcal{O}(nk \log^2 p)$	$\mathcal{O}(nk \log p)$

### 3 Preliminaries

In this section, we present the background needed for understanding the technical details of the rest of this paper. The main part of this section is to formally define the security of top- $\kappa$  protocols. We also present cryptographic tools used in our constructions in this paper.

*Notation.* Let us denote  $F(\alpha)$  for the multiplicity (also known as frequency) of an element  $\alpha$  in a multiset  $\mathbf{X}$  and  $F(\mathbf{X})$  for the collection of multiplicities for all elements in the multiset  $\mathbf{X}$ —here the multiplicity  $F(\alpha)$  of an element  $\alpha$  refers to how many times the element appears in  $\mathbf{X}$ .

For  $n \in \mathbb{N}$ ,  $[1, n]$  denotes the set  $\{1, \dots, n\}$ . If  $A$  is a probabilistic polynomial-time (PPT) machine, we use  $a \leftarrow A$  to denote  $A$  which produces output according to its internal randomness. In particular, if  $U$  is a set, then  $r \xleftarrow{\$} U$  is used to denote sampling from the uniform distribution on  $U$ . A function  $g : \mathbb{N} \rightarrow \mathbb{R}$  is negligible if for every positive polynomial  $\mu(\cdot)$  there exists an integer  $L$  such that  $g(\eta) < 1/\mu(\eta)$  for all  $\eta > L$ .

#### 3.1 Definitions

Informally, a double encryption is a pair of encryption schemes  $\mathcal{E} = (\text{KG}, \text{Enc}, \text{Dec})$  and  $\mathbf{E} = (G, E, D)$  such that  $\text{Enc}(\mathbf{E}(a)) = E(\text{Enc}(a))$ . We demand that an encryption scheme  $\mathbf{E}$  be *deterministic*. The reason is that we need to know a complete distribution of multisets while hiding their elements. See Appendix A for a formal definition of public-key cryptosystem and its standard security definition.

**Definition 1 (Double Encryption)** Let  $\mathcal{E} = (\text{KG}, \text{Enc}, \text{Dec})$  be a public-key encryption scheme defined as in Definition 4 with a pair of keys  $(pk, sk) \leftarrow \text{KG}(1^\lambda)$  and a message space (resp., ciphertext space)  $\mathbf{M}_{pk}$  (resp.,  $\mathbf{C}_{pk}$ ). A pair  $(\mathcal{E}, \mathbf{E})$  is called double encryption if there exists a triple of polynomial-time computable functions,  $\mathbf{E} = (G, E, D)$ , that satisfies the following properties:

- A probabilistic function  $G(1^\lambda)$  takes as input a security parameter  $\lambda$ , and outputs  $(s, s')$  such that for all  $s, s'$  and for any  $m \in \mathbf{M}_{pk}$ ,  $m = D_{s'}(E_s(m))$ . We demand that  $E$  and  $D$  be deterministic.
- *Commutativity:* For all  $pk, s$  and for all  $m \in \mathbf{M}_{pk}$ ,  $\text{Enc}_{pk}(E_s(m)) = E_s(\text{Enc}_{pk}(m))$  up to the randomness of  $\text{Enc}_{pk}(\cdot)$ .
- For all  $c \leftarrow \text{Enc}(m)$ ,  $E_s(m) = \text{Dec}_{sk}(E_s(c))$ .

We give an instantiation of a double encryption scheme in the following example.

**Example 1** Let  $p$  be a large prime of the form  $p = 2q + 1$ , where  $q$  is also prime. Let  $\mathbb{G}_q$  be a subgroup of  $\mathbb{Z}_p^\times$  of order  $q$  with a generator  $g$ . Then a standard CPA-secure ElGamal encryption  $\mathcal{E} = (\text{KG}, \text{Enc}, \text{Dec})$  is defined as follows: Selecting  $x \xleftarrow{\$} \mathbb{Z}_q$ ,  $\text{KG}(1^\lambda)$  outputs  $(pk, sk)$  where  $pk := (p, q, g, y = g^x \pmod{p})$  and  $sk := (p, q, g, x)$ . Given a message  $m \in \mathbb{G}_q$ , the encryption algorithm  $\text{Enc}$  outputs  $c = (g^r, m \cdot y^r)$  for a randomness  $r \xleftarrow{\$} \mathbb{Z}_q$ . Given an ElGamal ciphertext  $c = (u, v)$ , the decryption algorithm  $\text{Dec}$  computes  $v \cdot u^{-x}$  using the secret key  $x$ . Now  $\mathbf{E} = (G, E, D)$  is defined as follows:

- A probabilistic function  $G(1^\lambda)$  outputs  $(s, s') \in (\mathbb{Z}_q)^2$  such that  $ss' = 1 \pmod{q}$ .
- Given a value  $\alpha \in \mathbb{G}_q$ , a deterministic function  $E : \mathbb{G}_q \rightarrow \mathbb{G}_q$  is given by  $\alpha \mapsto \alpha^s \pmod{p}$ .
- A deterministic function  $D_{s'}(\beta)$  computes  $\beta^{s'} \pmod{p}$ .

Then,  $(\mathcal{E}, \mathbf{E})$  is a *double encryption*:

- For all values  $m \in \mathbb{G}_q$ ,  $m = (m^s)^{s'} \pmod{p}$ .
- For any message  $m \in \mathbb{G}_q$ , there exists a randomness  $r' = rs$  such that  $(g^{r'}, (m^s) \cdot y^{r'}) = ((g^r)^s, (my^r)^s)$ .
- For any ElGamal ciphertext  $c = (u, v) \in (\mathbb{G}_q)^2$ , where  $u = g^r$  and  $v = my^r$ ,  $m^s = v^s \cdot (u^s)^{-x} \pmod{p}$ .

We use a standard definition of shuffle given by Nguyen et al. [25]; see details of the definition in [25, Def. 3 & Def. 4] and its extend version. As we mentioned before, our shuffle algorithm takes as input a list of ciphertexts, and outputs a list of *permuted* and *doubly encrypted ciphertexts*. Since a double encryption scheme leads to change the plaintexts of input ciphertexts, we need to devise proving the correctness of the shuffle.

Now we define top- $\kappa$  protocol and give its algorithmic description. Throughout the paper, we use  $\Sigma_n$  to denote the set of all permutations on  $[1, n]$ . A private top- $\kappa$  protocol consists of five computable (in polynomial time) algorithms, (**Setup**, **DEncrypt**, **Shuffle**, **Aggregate**, **Reveal**), over a double encryption  $(\mathcal{E}, \mathbf{E})$ , which are explained as follows:

- $(pk, sk, s, s') \leftarrow \mathbf{Setup}(1^\lambda)$ : The setup algorithm takes as an input the security parameter  $\lambda$ , and outputs the public and secret parameters for doubly encrypting input ciphertexts, by invoking  $(pk, sk) \leftarrow \mathbf{KG}(1^\lambda)$  and  $(s, s') \leftarrow G(1^\lambda)$ .
- $(E_s(c_1), \dots, E_s(c_n)) \leftarrow \mathbf{DEncrypt}(pk, s, c_1, \dots, c_n)$ : The algorithm **DEncrypt** takes as input system parameters  $(pk, s)$  and a list of ciphertexts  $(c_1, \dots, c_n)$ , and produces a list of doubly encrypted ciphertexts  $(E_s(c_1), \dots, E_s(c_n))$ .
- $(E_s(c_{\pi(1)}), \dots, E_s(c_{\pi(n)})) \leftarrow \mathbf{Shuffle}(\pi, E_s(c_1), \dots, E_s(c_n))$ : The algorithm **Shuffle** chooses a random permutation  $\pi \in \Sigma_n$  and shuffles the doubly encrypted ciphertexts  $(E_s(c_1), \dots, E_s(c_n))$ , and then outputs the mixed list.
- $(E_s(\alpha_1), \dots, E_s(\alpha_\kappa)) \leftarrow \mathbf{Aggregate}(pk, sk, E_s(c_{\pi(1)}), \dots, E_s(c_{\pi(n)}))$ : The algorithm **Aggregate** takes as input a pair of keys  $(pk, sk)$  and a list of permuted, doubly encrypted ciphertexts, and for all  $i \in [1, n]$  computes  $\mathbf{Dec}_{sk}(E_s(c_{\pi(i)})) = E_s(\alpha_{\pi(i)})$ . Finally it computes  $F(E_s(\alpha_{\pi(i)}))$ ,  $i \in [1, n]$  and outputs only the elements whose multiplicity is greater than or equal than  $\kappa$ . Here  $\{E_s(\alpha_1), \dots, E_s(\alpha_\kappa)\} = \{E_s(\alpha_{\pi(i)}) \mid F(E_s(\alpha_{\pi(i)})) \geq \kappa\}$ .
- $(\alpha_1, \dots, \alpha_\kappa) \leftarrow \mathbf{Reveal}(pk, s', E_s(\alpha_1), \dots, E_s(\alpha_\kappa))$ : The algorithm **Reveal** outputs a set of the most frequent top- $\kappa$  elements by computing  $D_{s'}(E_s(\alpha_j))$  for all  $j \in [1, \kappa]$ .

In the next section, we will define the meaning of *secure* top- $\kappa$  protocol. Then we describe cryptographic building blocks for constructing a secure top- $\kappa$  protocol under proper cryptographic assumptions.

### 3.2 Security Definition

*Ideal Functionality.* Firstly, we define the ideal functionality  $\mathcal{F}_{\text{top}\kappa}$  for top- $\kappa$  protocol as follows:

**Definition 2** *There are a set of  $n$  users,  $U = \{u_i\}_{i=1}^n$ , a trusted party  $\mathcal{T}$ , and an ideal adversary  $\mathcal{S}$  controlling a set of corrupted users  $\Upsilon_t = \{u_{i_j}\}_{j=1}^t$  for some  $t \in [0, n-1]$ . Let  $\mathbf{X}_i = \{\alpha_{i,j}\}_{j=1}^{k_i}$  be a multiset of user  $u_{i \in [1, n]}$ .<sup>6</sup>*

(1) *Each user  $u_i$  sends  $\mathbf{X}_i$  to  $\mathcal{T}$ .*

(2)  *$\mathcal{T}$  computes the following functionality, and returns the output  $Z_l$  to each  $u_{l \in [1, n]}$ :*

$$Z_l = \left\{ \alpha_{i,j} \in \bigcup_{i \in [1, n]} \mathbf{X}_i \mid F(\alpha_{i,j}) \geq \kappa \right\}.$$

<sup>6</sup> In fact, in our setting  $k_i = k_{i'}$  for all  $i, i' \in [1, n]$ .

*Data Privacy.* Informally, *data privacy* requires that no user, or coalition of users, should learn anything about honest users' inputs except what can be trivially derived from the output itself. We can easily derive the formal definition for data privacy in top- $\kappa$  protocols following the standard privacy definition of existing protocols in the literature; an excellent reference on that is Goldreich's textbook in [13]. More specifically, we use Definition 7.5.1 (resp., Definition 7.5.3) in [13] for the semi-honest model (resp. the malicious model). Roughly speaking, this is formalized by considering an ideal world where  $\mathcal{T}$  receives the inputs of the users and outputs the result of the defined functionality. We demand that in the real world application of the protocol—that is, one without the  $\mathcal{T}$ —no user should learn more information than in the ideal world.

*User Privacy.* The remaining part to conclude our definitions is user privacy. Let  $Z = \{\alpha_1, \dots, \alpha_\kappa\}$  be an output of a top- $\kappa$  protocol. Roughly speaking, no user or coalition of users should gain a non-negligible advantage in distinguishing, for all  $\alpha \in Z$ , an honest user  $u_i$  such that  $\alpha \in X_i$ .

**Definition 3 (User Privacy)** Let  $\Pi_{\kappa, \mathcal{E}, \mathbf{E}}$  be a top- $\kappa$  protocol defined as in Section 3.1 over a double encryption scheme  $(\mathcal{E}, \mathbf{E})$  and  $\mathcal{A} = (\mathcal{A}_1, \mathcal{A}_2)$  be an adversary.

*Experiment*  $\text{Exp}_{\mathcal{A}}^{\text{top-}\kappa}(\Pi_{\kappa, \mathcal{E}, \mathbf{E}}, \lambda)$   
 $(pk, sk, s, s') \leftarrow \text{Setup}(\lambda);$   
 $(\text{state}, \mathcal{Y}_t, m_1, \dots, m_{n-t}) \leftarrow \mathcal{A}_1(pk, n, t)$  where  $\mathcal{Y}_t$  is a set of corrupted  $t$  users;  
 $\sigma \xleftarrow{\$} \Sigma_n$  and assign  $m_{\sigma(i)}$  to each honest  $i$ -th user  $u_i \notin \mathcal{Y}_t$ ;  
 $(\alpha_1, \dots, \alpha_\kappa) \leftarrow \Pi_{\kappa, \mathcal{E}, \mathbf{E}}$ , where  $\mathcal{A}_1$  interacts with the  $n - t$  honest users;  
 $(i, j) \leftarrow \mathcal{A}_2(pk, m_1, \dots, m_{n-t}, \text{state});$

We define the advantage of an adversary  $\mathcal{A}$ , running in probabilistic polynomial time, as:

$$\text{Adv}_{\mathcal{A}}^{\text{top-}\kappa}(\Pi_{\kappa, \mathcal{E}, \mathbf{E}}, \lambda) = \left| \Pr[\sigma(i) = j] - \frac{1}{n-t} \right|.$$

A top- $\kappa$  protocol is user-private if the advantage  $\text{Adv}_{\mathcal{A}}^{\text{top-}\kappa}(\Pi_{\kappa, \mathcal{E}, \mathbf{E}}, \lambda)$  is negligible in the security parameter  $\lambda$ .

In conclusion, a top- $\kappa$  protocol  $\Pi_{\kappa, \mathcal{E}, \mathbf{E}}$  over a double encryption  $(\mathcal{E}, \mathbf{E})$  is said to be *secure* if  $\Pi_{\kappa, \mathcal{E}, \mathbf{E}}$  is data-private and user-private.

### 3.3 Cryptographic Assumptions and Tools

In the following we outline the cryptographic tools and assumptions we use in our protocol. Let  $\mathcal{G}$  be a finite cyclic group of prime order  $q$ , and let  $g \in \mathcal{G}$  be a generator. Given  $h \in \mathcal{G}$ , the discrete logarithm problem requires us to compute  $x \in \mathbb{Z}_q$  such that  $g^x = h$ . We denote this (unique)  $x$  by  $\log_g h$ . In particular groups  $\mathcal{G}$  and for  $q$  large, it is assumed hard to compute  $x$ , which is said to be the Discrete Logarithm (DL) assumption.

A stronger assumption is the Decisional Diffie-Hellman (DDH) assumption. Here, given  $\mathcal{G}$ , a generator  $g$  of  $\mathcal{G}$ , and three elements  $a, b, c \in \mathcal{G}$ , we are asked (informally) to decide whether there exist  $x, y$  such that  $a = g^x, b = g^y$ , and  $c = g^{xy}$ . More formally, the DDH assumption states that the following two distributions are computationally indistinguishable:  $\{\mathcal{G}, g, g^x, g^y, g^{xy}\}$  and  $\{\mathcal{G}, g, g^x, g^y, g^z\}$  where  $x, y, z \xleftarrow{\$} \mathbb{Z}_q$ .

We extensively use ElGamal encryption defined in Example 1. This ElGamal encryption scheme is defined and known to be secure against CPA attack in a DDH group  $\mathbb{G}$ ; see Appendix A for the CPA security. In addition, we need an efficient scheme which works as follows: When each user holds a shared secret key  $s_i$  such that  $s = \prod_{i=1}^n s_i$ , the scheme allows each user to have a share  $s'_i$  satisfying  $s' = \prod_{i=1}^n s'_i$  and  $s' = s^{-1} \pmod{q}$  for a public modulus  $q$ . Indeed, we may realize the scheme by techniques studied by Algesheimer et al. [3, §5].

## 4 Our Construction

In this section, we describe our construction for computing the top- $\kappa$  elements privately. We begin by considering a basic setting of  $n$  users, denoted by  $u_1, \dots, u_n$ . Let  $\mathbf{X}_i = \{\alpha_{i,1}, \dots, \alpha_{i,k}\}$  for all  $i \in [1, n]$ . Each user  $u_i$  has its private multiset  $\mathbf{X}_i$ , and the users wish to jointly compute  $\{\alpha \in \bigcup_{i=1}^n \mathbf{X}_i \mid F(\alpha) \geq \kappa\}$ . For simplicity, assume that all elements are in the proper message domain  $\mathbf{M}_{pk}$  of an ElGamal encryption scheme, e.g., a finite cyclic subgroup  $\mathbb{G}_q$  of  $\mathbb{Z}_p^\times$  in which the DDH assumption holds. For a multiset  $\mathbf{X} = \{\alpha_1, \dots, \alpha_k\}$ , we denote  $\mathbf{X}^s$  as  $\{\alpha_1^s, \dots, \alpha_k^s\}$  for some  $s \in \mathbb{Z}_q$ . With such notation in mind, we proceed to describe our construction.

### 4.1 Description

Let  $\lambda$  be a security parameter that will determine the size of groups. Let  $p$  be a  $\lambda$ -bit prime such that for some prime  $q$ ,  $p = 2q + 1$ . Let  $\mathbb{G}_q$  be a finite cyclic subgroup of  $\mathbb{Z}_p^\times$  whose order is  $q$ , and let  $g$  be a generator of  $\mathbb{G}_q$ . Our construction follows:

**Setup**( $1^\lambda$ ) Each user agrees to a threshold ElGamal encryption  $\mathcal{E}$  with a public/private key pair  $(pk, sk)$ , which are computed as follows. Define  $\text{params} := (p, q, g, \mathbb{G}_q)$ . Each user selects a value  $x_i \xleftarrow{\$} \mathbb{Z}_q$ , computes  $y_i = g^{x_i}$ , and sets  $sk = (\text{params}, x_i)$ ; the public key is then given by  $pk := (\text{params}, y = \prod_{i=1}^n y_i = g^{\sum_{i \in [1, n]} x_i} \pmod{p})$ .<sup>7</sup> In addition, all users are distributed a share  $(s_i, s'_i)$  such that  $s = \prod_{i=1}^n s_i$ ,  $s' = \prod_{i=1}^n s'_i$ , and  $s \cdot s' = 1 \pmod{q}$ .

**DEncrypt** Let  $I = \{1, \dots, n\}$  be a set of indices, and let the power function  $E_{s_i}(\alpha) = \alpha^{s_i} \pmod{p}$  which is deterministic.

1. Every user  $u_i$  encrypts his multiset  $\mathbf{X}_i$  as follows:

$$\text{Enc}_{pk}(\mathbf{X}_i) = \{\text{Enc}_{pk}(\alpha_{i,1}), \dots, \text{Enc}_{pk}(\alpha_{i,k})\}$$

where  $\text{Enc}_{pk}(\alpha_{i,j}) = (g^{r_{i,j}}, \alpha_{i,j} \cdot y^{r_{i,j}})$  for some randomizer  $r_{i,j} \in \mathbb{Z}_q$ , and sends  $\text{Enc}_{pk}(\mathbf{X}_i)$  to  $u_1$ .

2. User  $u_1$  computes  $\{E_{s_1}(\text{Enc}_{pk}(\mathbf{X}_1)), \dots, E_{s_1}(\text{Enc}_{pk}(\mathbf{X}_n))\}$ , which is denoted by  $Y_0$ .

**Shuffle & DEncrypt** For  $i \in [1, n]$ ,  $u_i$  receives vector  $Y_{i-1}$  and computes a permuted, doubly encrypted version  $Y_i$  as follows:

1.  $u_{i \neq 1}$  computes

$$\begin{aligned} E_{s_i}(Y_{i-1}) &= \{c_1, \dots, c_{nk}\} \\ &= \{E_{s_i}(E_{s_{i-1}}(\dots E_{s_1}(\alpha_{\pi_{i-1}(1)}) \dots)), \dots, E_{s_i}(E_{s_{i-1}}(\dots E_{s_1}(\alpha_{\pi_{i-1}(nk)}) \dots))\}. \end{aligned}$$

More precisely, here  $\alpha_{\pi_{i-1}(\ell)} = \alpha_{\pi_{i-1} \circ \dots \circ \pi_1(\ell)}$  for all  $\ell \in [1, nk]$ .

2.  $u_i$  chooses a random permutation  $\pi_i \in \Sigma_{nk}$ , and applies  $\pi_i$  to the list of  $c_{\ell \in [1, nk]}$  computed above; denote the result by  $Y_i$ .
3.  $u_i$  sends  $Y_i$  to  $u_{i+1}$ ; the last user  $u_n$  sends  $Y_n$  to all users.

**Aggregate** Let  $\mathbf{U} = \bigcup_{i=1}^n \mathbf{X}_i$ . Every user has  $E_s(\text{Enc}_{pk}(\mathbf{U}))$ .

1. Every user participates in a group decryption and obtains

$$E_s(\mathbf{U}) = \{E_s(\alpha_{\pi(1)}), \dots, E_s(\alpha_{\pi(nk)})\}$$

where  $\pi = \pi_n \circ \dots \circ \pi_1$ .

2. Every user computes  $Z = \{E_s(\alpha) \in E_s(\mathbf{U}) \mid F(E_s(\alpha)) \geq \kappa\}$ .

**Reveal**

1. For every  $E_s(\alpha) \in Z$ , user  $u_i$  sends its share of  $D_{s'_i}(E_s(\alpha))$  to  $u_i'$ .
2. After receiving all the shares, every user  $u_i$  computes  $\alpha = D_{s'}(E_s(\alpha))$ , thereby recovering the top- $\kappa$ ,  $\{\alpha \in \mathbf{U} \mid F(\alpha) \geq \kappa\}$ .

<sup>7</sup> Notice that in threshold decryption schemes, users generally produce shares of the decrypted element, and during the operation of the schemes if one user sends a uniformly generated share instead of a valid one the decrypted element is uniform. Also, if the decrypted element is uniform, the resulting decryption reveals no information to the users.

**Efficiency.** The advantage of the above protocol is multifold. First, compared to Kissner and Song’s protocol [22], our protocol provides the functionality of finding a threshold value and computing the “over threshold” at the same computation and communication cost—whereas they incur different and higher costs in [22]. Second, compared to the top- $\kappa$  protocol described in [6], our protocol has a much better computational complexity. See details in Section 5.

In order to present a fair comparison between our proposed protocol and Applebaum et al.’s protocol [4] we devise our protocol for a semi-decentralized model in the next section. The other purpose of our modification is to reduce the round complexity to a constant.

## 4.2 Semi-Decentralized Construction

The most crucial drawback of our proposed protocol in the previous subsection is that it has  $\mathcal{O}(n)$  round complexity. To avoid this problem, Applebaum et al. introduced two semi-honest users: a *proxy* which shuffles a list of input ciphertexts, and a *database* which aggregates top- $\kappa$  elements. Applying the same technique to our protocol, we obtain a constant-round top- $\kappa$  protocol. We briefly summarize the major changes in the following:

- Assume that there are  $n_1$  proxies and  $n_2$  databases described as in [4].
- Each database engages in setting up a threshold ElGamal encryption and publishes a public key. Instead of users, all proxies are distributed secret shares  $(s_l, s'_l)_{l \in [1, n_1]}$  as above.
- Each user computes a list of ElGamal ciphertexts and sends it to a proxy.
- Each proxy runs DEncrypt and Shuffle algorithms, and returns the result to all databases.
- Databases perform a group decryption, find a list of encrypted top- $\kappa$  elements, and output the list.
- Finally all proxies decrypt the encrypted top- $\kappa$  list and return the top- $\kappa$  to all users.

Compared to [4], our protocol does not require OT operations between users and proxies, nor an extra encryption scheme. Recall that Applebaum et al.’s protocol requires ElGamal encryption and Goldwasser-Micali (GM) encryption: ElGamal encryption is used to encrypt elements in multisets and GM encryption is used to encrypt their multiplicity.

## 5 Evaluation

In this section, we provide security and complexity analyses of our construction. We first prove the correctness of our construction, so that every semi-honest user participating in the protocol learns the joint distribution of users’ multisets and the number of times each of the elements appears, with overwhelming probability. Furthermore, we show that our construction is secure and data of honest users are maintained private as long as the size of the coalition among semi-honest users is less than  $n$ .

### 5.1 Security Analysis

**Theorem 1 (Correctness)** *In the private top- $\kappa$  protocol in Section 4.1, every honest user learns the joint set distribution of all users’ private inputs, i.e., each element  $E_s(\alpha)$  such that  $\alpha \in \bigcup_{i=1}^n \mathbf{X}_i$  and the number of times it appears, with overwhelming probability.*

*Proof.* Each player learns a randomly permuted joint multiset  $E_s(\mathbf{U}) = \{E_s(\alpha_{\pi(1)}), \dots, E_s(\alpha_{\pi(nk)})\}$ . We know that  $|\mathbf{U}^s| = nk$ . Since  $\pi$  is a permutation, for each  $E_s(\alpha_{\pi(\ell)})$  and for all  $\ell \in [1, nk]$ , there exist a pair of the unique index  $\ell^*$  such that

$$\begin{aligned} \ell^* &= \pi^{-1}(\ell) \\ &= \pi_n^{-1}(\ell) \circ \dots \circ \pi_1^{-1}(\ell). \end{aligned}$$

Namely,  $E_s(\alpha_{\pi(\ell)})$  is a unique blinded version of  $\alpha_{\ell^*} \in \bigcup_{i=1}^n \mathbf{X}_i$ . Moreover, for all  $\ell, \ell^* \in [1, nk]$ ,  $\alpha_\ell = \alpha_{\ell^*}$  if and only if  $E_s(\alpha_\ell) = E_s(\alpha_{\ell^*})$  with overwhelming probability.  $\square$

**Corollary 1** *In the private top- $\kappa$  protocol in Section 4.1, every honest user learns the top- $\kappa$  in the union of private multisets with overwhelming probability.*

Now we show that our protocol satisfies the privacy requirements in the semi-honest model. Let  $\mathcal{T}$  be a trusted party in the ideal world which receives the private input multiset  $\mathbf{X}_i$  of size  $k$  from user  $u_i$  for  $i \in [1, n]$ , and then returns to every user the joint multiset distribution  $\{F(\alpha)\}$  for all  $\alpha \in \bigcup_{i=1}^n \mathbf{X}_i$ .

**Theorem 2 (Data Privacy)** *Assume that the threshold ElGamal encryption  $\mathcal{E} = (\text{KG}, \text{Enc}, \text{Dec})$  is secure against CPA. In the private top- $\kappa$  protocol in Section 4.1, any coalition of less than  $n$  semi-honest users learn no more information than would be given by using the same private inputs in the ideal-world model with  $\mathcal{T}$ .*

*Proof.* We assume that the ElGamal encryption scheme is CPA-secure, and so each user learns only

$$\begin{aligned} & \text{Enc}_{pk}(\mathbf{X}_1), \dots, \text{Enc}_{pk}(\mathbf{X}_n); \\ & E_{s_1}(\text{Enc}_{pk}(\mathbf{X}_1)), \dots, E_{s_1}(\text{Enc}_{pk}(\mathbf{X}_1)), \dots, E_{s_{i-1}}(\text{Enc}_{pk}(\mathbf{X}_1)), \dots, E_{s_{i-1}}(\text{Enc}_{pk}(\mathbf{X}_n)); \\ & \vdots \\ & E_{s_{i-1}}(\dots E_{s_1}(\text{Enc}_{pk}(\mathbf{X}_1)) \dots), \dots, E_{s_{i-1}}(\dots E_{s_1}(\text{Enc}_{pk}(\mathbf{X}_n)) \dots) \end{aligned}$$

during an execution. At the end of the protocol all users further know  $E_s(\text{Enc}_{pk}(\mathbf{U}))$  where  $\mathbf{U} = \bigcup_{i=1}^n \mathbf{X}_i$ , and for some  $\gamma_{\ell \in [1, nk]} \in \mathbb{Z}_q$

$$\begin{aligned} E_s(\text{Enc}_{pk}(\mathbf{U})) &= \{E_s(\text{Enc}_{pk}(\mathbf{X}_1)), \dots, E_s(\text{Enc}_{pk}(\mathbf{X}_n))\} \\ &= (g^{\gamma^1}, (\alpha_{\pi(1)})^s \cdot y^{\gamma^1}), \dots, (g^{\gamma^{nk}}, (\alpha_{\pi(nk)})^s \cdot y^{\gamma^{nk}}). \end{aligned}$$

Note that  $\pi$  is a composition of random permutations and is unknown to all users, as the maximum coalition size is smaller than  $n$ . That is, if there exists at least an honest user, then a composition of random permutations  $\pi = \pi_n \circ \dots \circ \pi_1$  is a random permutation because at least a permutation  $\pi_{i \in [1, n]}$  is secure. What is more, note that  $s$  is uniformly distributed and unknown to all users for the same reason. As  $s$  is uniformly distributed for any user inputs and  $\pi$  is random, no user or coalition can learn more than a set of re-randomized ElGamal encryptions. As  $s$  is uniformly distributed, a group decryption of ElGamal encryptions reveals no more than

$$\{E_s(\alpha_\ell)\}_{\ell \in [1, nk]} = E_s\left(\bigcup_{i=1}^n \mathbf{X}_i\right) = E_s(\mathbf{U}).$$

We know the fact that  $F(\alpha) = F(\alpha^s)$  for two multisets  $\mathbf{X}$  and  $E_s(\mathbf{X}) \in (\mathbb{G}_q)^k$ , for all  $s \in \mathbb{Z}_q$  and for all  $\alpha \in \mathbf{X}$ . Hence we see that

$$F(E_s(\mathbf{U})) = F\left(E_s\left(\bigcup_{i=1}^n \mathbf{X}_i\right)\right) = F\left(\bigcup_{i=1}^n \mathbf{X}_i\right) = F(\mathbf{U}),$$

which can be derived from the output that would be returned by  $\mathcal{T}$  in the ideal-world model.  $\square$

Finally we show that our scheme guarantees user privacy. Intuitively, this proof ensures that even when  $n - 1$  users collude with each other in an arbitrary way they cannot identify their own elements from the joint multiset.

**Theorem 3 (User Privacy)** *Assume that the threshold ElGamal encryption  $\text{Enc}$  is CPA-secure. The private top- $\kappa$  protocol in Section 4.1 is user-private against any coalition of less than  $n$  semi-honest users.*

*Proof.* Assume that there is at least an honest user in the system, and that the threshold ElGamal encryption  $\mathcal{E} = (\text{KG}, \text{Enc}, \text{Dec})$  is CPA-secure. After performing DEncrypt and Shuffle algorithms, every user obtains a collection of ElGamal encryptions  $\{c_1, \dots, c_{nk}\}$ . By the second assumption, the adversary cannot learn any further information except that which encryptions have been sent from which users. Running these algorithms, each user should raise the power of the received encryptions with his shared secret  $s_i$ . Namely, each user holds the modified list of the encryptions,

$$\{E_{s_i}(c_1), E_{s_i}(c_2), \dots, E_{s_i}(c_{nk})\}.$$

Next the user should apply his private permutation  $\pi_i$  to the list so that she gets changed into

$$\{E_{s_i}(c_{\pi(1)}), E_{s_i}(c_{\pi(2)}), \dots, E_{s_i}(c_{\pi(nk)})\}.$$

At the end of running these algorithms, all users get a permuted and doubly encrypted list

$$\{E_s(c_{\pi(1)}), E_s(c_{\pi(2)}), \dots, E_s(c_{\pi(nk)})\}$$

where the permutation  $\pi = \pi_n \circ \dots \circ \pi_1$  and  $s = \prod_{i=1}^n s_i$ . As there exists at least an honest user, even though all  $n - 1$  users collude with each other  $s$  is uniformly distributed and unknown to all users and  $\pi$  is a random permutation. This completes the proof of the claim.  $\square$

**Theorem 4** *Assuming that the threshold ElGamal encryption is CPA-secure and the DL assumption holds, the proposed top- $\kappa$  protocol is secure in the semi-honest model.*

*Proof.* We complete the proof of security by Theorem 2 and Theorem 3.  $\square$

## 5.2 Efficiency Analysis

The private top- $\kappa$  protocol has not yet been implemented, but we give a detailed analysis of the running time and space requirements as follows. We base our protocol on ElGamal encryption and the power function with primes  $|p| = 1024$ ,  $|q| = 160$ . To measure users' computational loads, we count the number of exponentiations using a 1024-bit modulus.

**Table 2.** Complexity Analysis

	Comp. Cpx (expo.)	Comm. Cpx (bits)	Rounds Cpx
Setup	$n$	$n \log p$	1
DEncrypt & Shuffle	$4nk + 2n^2k$	$2(n-1)k \log p + 2n^2k \log p$	$n$
Aggregate	$n^2k$	$2n^2k \log p$	1
Reveal	$n\kappa$	$n\kappa \log p$	$n - 1$

In Table 2 we show a summary of the complexity result for our proposed protocol. The total computational complexity is dominated by DEncrypt and Shuffle algorithms. Putting the computational complexities together shows that total computation complexity is  $\mathcal{O}(n^2k)$  in  $\mathcal{O}(n)$  rounds. The proposed protocol has  $\mathcal{O}(n^2k \log p)$  bits in total. We compare our protocol with the most relevant ones in the following. Note that Applebaum et al.'s scheme [4] runs in the semi-decentralized model. Since it is not fair to directly compare our protocol with Applebaum et al.'s protocol, we just present the computational complexity.

**Comparison.** We consider three existing protocols: one based on Kissner and Song (KS) protocol [22], Burkhart and Dimitropoulos (BD) protocol [6], and Applebaum et al. protocol.

**Based on KS protocol** We first compare our scheme with a top- $\kappa$  protocol that can be constructed by using KS protocol. As mentioned above, since it does not provide a way for finding a threshold value  $\tau$ , we do not know computational and communication complexity in computing  $\tau$ . Assuming  $\tau$  is given, their protocol has  $\mathcal{O}(n^2k)$  computation complexity in  $\mathcal{O}(n)$  rounds.

**BD protocol** In turn, we give a comparison with BD protocol. To our knowledge, it is the only fully decentralized top- $\kappa$  protocol that does not use Yao’s garbled circuit evaluation. Their protocol utilizes two special-purpose sub-protocols—**equality** and **lessthan** (see [9,26]), but in [7] as the authors pointed out, comparison of two shared secrets is very expensive and computational intensive. Thus, they use a computationally efficient version of the basic sub-protocols as follows: **equality** requires  $\log p$  rounds and **lessthan** requires  $(2 \log p + 10)$  rounds. Their protocol consists of two key ingredients as follows:

- Finding a correct threshold value  $\tau$ : Requires  $(\log k(2 \log p + 10) + \log p + 2 \log p + 10) nk$  rounds.
- Resolving collisions: Requires  $\frac{n(n-1)}{2} \log p + 2(n-1) \log p + 10(n-1)$  rounds.

Note that BD protocol also should know  $\tau$  as in KS protocol. Hence, the total round complexity is  $\mathcal{O}(n(n+k \log k) \log p)$  where we set the size of hash tables by  $\log k$  and the maximum size of union  $U$  by  $nk$ .

With a few simple calculations, we can see that their protocol incurs  $4 \left( \frac{n(n-1)}{2} k + k(n-1) \right)$  multiplications in  $\mathbb{Z}_p^\times$ .

**Applebaum et al. protocol** Let us use  $\mathcal{O}_p(\cdot)$  to denote complexity using modulus prime  $p$  and  $\mathcal{O}_N(\cdot)$  complexity using modulus composite  $N$ . Assume all elements are integers less than  $p$  and the maximum multiplicity is less than  $\log \log p$ .

Their major computation-intensive parts are as follows:

- Interactive computation between Users and Proxy: First, users should run a protocol for oblivious evaluation of pseudorandom function by communicating with proxies, and then encrypt the result with ElGamal encryption. This requires  $n(k(2 \log p + 2) + 2k)$  exponentiations over  $\mathbb{Z}_p^\times$ . Further, users should encrypt the multiplicity of each element with GM encryption, which requires  $nk \log \log p$  multiplications over  $\mathbb{Z}_N^\times$ . Finally each user doubly encrypts their elements using ElGamal encryption. This requires  $2nk$  exponentiations over  $\mathbb{Z}_p^\times$ .
- Aggregation by Database: The most computationally-intensive part is ElGamal and GM decryption. Since database receives two types of ElGamal ciphertexts, it performs  $2nk$  exponentiations over  $\mathbb{Z}_p^\times$ . GM decryption requires  $2nk \log \log p$  exponentiations over  $\mathbb{Z}_N^\times$ .

Therefore, the total computational complexity is  $\mathcal{O}_p(nk \log p) + \mathcal{O}_N(nk \log \log p)$  in terms of the number of exponentiations.

## 6 Handling Malicious Users

Our protocol is secure in the semi-honest model. Although, we describe possible modifications to our top- $\kappa$  protocol in order to provide security in the malicious adversary model and leave a realization of such protocol as a future work.

A general technique to prevent malicious behavior is to use zero-knowledge proofs. For example, in KS protocol whenever each user sends a public ciphertext  $\text{Enc}_{pk}(m)$ , it should prove that it knows the corresponding plaintext  $m$ . First, there exist efficient zero-knowledge proofs of knowledge of discrete logarithms (PODL) in our setting [8]. In **Setup** algorithm each user proves that given a public-key share  $y_i$  it knows the corresponding secret key  $x_i$  using PODL. Each user proves in a zero-knowledge manner that given an ElGamal ciphertext it knows the corresponding message by zero-knowledge proof of knowledge of randomizers, which is also a standard technique. Furthermore, it seems that we can prove the correctness of shuffling a list of ElGamal ciphertexts by modifying Groth’s technique [15]. The challenging part is that in our protocol **DEncrypt** algorithm forces the plaintexts in the input ciphertexts to be changed. Hence, we should be careful in designing zero-knowledge proofs for correctness of a shuffle.

## 7 Conclusion

In this paper we have looked at the problem of finding the top- $\kappa$  element securely, and formally defined what it means for a protocol to be a secure top- $\kappa$  protocol. We examined the existing constructions in the literature and compared them with our proposed protocol. We further

developed two protocols, with varying operation overhead, analyzed their security, and demonstrated their practicality and superiority (computational and security-wise) to other work in the literature.

In the near future, we will investigate two directions. First, as our constructions' security is proven in the semi-honest model—which is rationalized by the application domain, we will investigate constructions that are provably secure in the malicious model, and their potential applications. Second, as the shuffling algorithm in our current construction requires sending messages among players in a relay manner, we will consider the practical and security aspects of a construction that relies on sending such messages in a broadcast manner.

## References

1. B. Adida. *Advances in cryptographic voting systems*. PhD thesis, Massachusetts Institute of Technology, 2006. 3
2. G. Aggarwal, N. Mishra, and B. Pinkas. Secure computation of the  $k^{\text{th}}$ -ranked element. In C. Cachin and J. Camenisch, editors, *Advances in Cryptology-EuroCrypt*, LNCS 3027, pages 40–55, 2004. 2
3. J. Algesheimer, J. Camenisch, and V. Shoup. Efficient computation modulo a shared secret with application to the generation of shared safe-prime products. In M. Yung, editor, *Advances in Cryptology-Crypto*, LNCS 2442, pages 417–432, 2002. 7
4. B. Applebaum, H. Ringberg, M. Freedman, M. Caesar, and J. Rexford. Collaborative, privacy-preserving sata aggregation at scale. In M. Atallah and N. Hopper, editors, *Privacy Enhancing Technologies*, LNCS 6205, pages 56–74, 2010. 1, 3, 4, 5, 9, 11
5. M. Bellare and P. Rogaway. Random oracles are practical: a paradigm for designing efficient protocols. In D. Denning, R. Pyle, R. Ganesan, R. Sandhu, and V. Ashby, editors, *ACM Conference on Computer and Communications Security*, pages 62–73, 1993. 4
6. M. Burkhart and X. Dimitropoulos. Fast privacy-preserving top- $k$  queries using secret sharing. In *International Conference on Computer Communications and Networks (ICCCN)*, 2010. 3, 4, 5, 9, 11
7. M. Burkhart, M. Strasser, D. Many, and X. Dimitropoulos. SEPIA: Privacy-preserving aggregation of multi-domain network events and statistics. In *USENIX Security Symposium*, pages 223–240, 2010. 12
8. J. Camenisch. Proof systems for general statements about discrete logarithms. Technical Report TR 260, Dept. of Computer Science, ETH Zurich, 1997. 12
9. I. Damgård, M. Fitzi, E. Kiltz, J. B. Nielsen, and T. Toft. Unconditionally secure constant-rounds multi-party computation for equality, comparison, bits and exponentiation. In S. Halevi and T. Rabin, editors, *TCC*, LNCS 3876, pages 285–304, 2006. 12
10. Y. Desmedt and K. Kurosawa. How to break a practical mix and design a new one. In B. Preneel, editor, *Advances in Cryptology-EuroCrypt*, LNCS 1807, pages 557–572, 2000. 1, 3
11. T. El Gamal. A public key cryptosystem and a signature scheme based on discrete logarithms. In G. R. Blakely and D. Chaum, editors, *Advances in Cryptology-Crypto*, LNCS 196, pages 10–18, 1984. 4
12. J. Furukawa and K. Sako. An efficient scheme for proving a shuffle. In J. Kilian, editor, *Advances in Cryptology-Crypto*, LNCS 2139, pages 368–387, 2001. 3
13. O. Goldreich. *The foundations of cryptography*, volume 2. Cambridge University Press, 2004. 7
14. S. Goldwasser and S. Micali. Probabilistic encryption. *J. Comput. Syst. Sci.*, 28(2):270–299, 1984. 5, 15
15. J. Groth. A verifiable secret shuffle of homomorphic encryptions. In Y. Desmedt, editor, *Public Key Cryptography*, LNCS 2567, pages 145–160, 2003. 12
16. J. Groth. A verifiable secret shuffle of homomorphic encryptions. *Journal of Cryptology*, 23(4):546–579, 2010. 3
17. J. Groth and S. Lu. Verifiable shuffle of large size ciphertexts. In T. Okamoto and X. Wang, editors, *Public Key Cryptography*, LNCS 4450, pages 377–392, 2007. 3
18. W. He, X. Liu, H. Nguyen, K. Nahrstedt, and T. F. Abdelzaher. PDA: Privacy-preserving data aggregation in wireless sensor networks. In *INFOCOM*, pages 2045–2053. IEEE, 2007. 1
19. J. Hong, J. W. Kim, J. Kim, K. Park, and J. H. Cheon. Constant-round privacy preserving multiset union. In *Cryptology ePrint Archive*, 2011/138, 2011. 2
20. Y. Ishai, J. Kilian, K. Nissim, and E. Petrank. Extending oblivious transfers efficiently. In D. Boneh, editor, *Advances in Cryptology-Crypto*, LNCS 2729, pages 145–161, 2003. 4

21. M. Jakobsson. A practical mix. In K. Nyberg, editor, *Advances in Cryptology-EuroCrypt*, LNCS 1403, pages 448–461, 1998. 1, 3
22. L. Kissner and D. Song. Privacy-preserving set operations. In V. Shoup, editor, *Advances in Cryptology-Crypto*, LNCS 3621, pages 241–257, 2005. 2, 4, 9, 11
23. M. Naor and B. Pinkas. Oblivious transfer with adaptive queries. In M. Wiener, editor, *Advances in Cryptology-Crypto*, LNCS 1666, pages 573–590, 1999. 4
24. C. Neff. A verifiable secret shuffle and its application to e-voting. In *ACM Conference on Computer and Communications Security*, pages 116–125, 2001. 1, 3
25. L. Nguyen, R. Safavi-Naini, and K. Kurosawa. Verifiable shuffles: A formal model and a Paillier-based efficient construction with provable security. In M. Jakobsson, M. Yung, and J. Zhou, editors, *ACNS*, LNCS 3089, pages 61–75, 2004. 3, 6
26. T. Nishide and K. Ohta. Multiparty computation for interval, equality, and comparison without bit-decomposition protocol. In T. Okamoto and X. Wang, editors, *Public Key Cryptography*, LNCS 4450, pages 343–360, 2007. 12
27. P. Paillier. Public-key cryptosystems based on composite degree residuosity classes. In J. Stern, editor, *Advances in Cryptology-EuroCrypt*, LNCS 1592, pages 223–238, 1999. 3
28. Y. Sang and H. Shen. Efficient and secure protocols for privacy-preserving set operations. *ACM Transactions on Information and System Security (TISSEC)*, 13(1):9:1–9:35, 2009. 2
29. J. Vaidya and C. Clifton. Privacy-preserving top- $k$  queries. In *ICDE*, pages 545–546, 2005. 2
30. Vijay S. Iyengar. Transforming data to satisfy privacy constraints. In *KDD*, pages 279–288, 2002. 1
31. L. Xiong, S. Chitti, and L. Liu. Top $k$  queries across multiple private databases. In *International Conference on Distributed Computing Systems (ICDCS)*, pages 145–154, 2005. 2
32. A. Yao. Protocols for secure computations. In *FOCS*, pages 160–164, 1982. 4
33. R. Zhang, J. Shi, Y. Liu, and Y. Zhang. Verifiable fine-grained top- $k$  queries in tiered sensor networks. In *INFOCOM*, pages 2633–2641, 2010. 2
34. R. Zhang, Y. Zhang, and C. Zhang. Secure top- $k$  query processing via untrusted location-based service providers. In *INFOCOM*, pages 1170–1178, 2012. 2

## A Basic Definitions

We first give a formal definition of a public key cryptosystem and then its standard security definition.

We shall write

$$\Pr[x_1 \stackrel{\$}{\leftarrow} X_1, x_2 \stackrel{\$}{\leftarrow} X_2(x_1), \dots, x_n \stackrel{\$}{\leftarrow} X_n(x_1, \dots, x_{n-1}) : \varphi(x_1, \dots, x_n)]$$

to denote the probability that when  $x_1$  is drawn from a certain distribution  $X_1$ , and  $x_2$  is drawn from a certain distribution  $X_2(x_1)$ , possibly depending on the particular choice of  $x_1$ , and so on, all the way to  $x_n$ , the predicate  $\varphi(x_1, \dots, x_n)$  is true.

**Definition 4** A public-key cryptosystem  $\mathcal{E}$  is a 3-tuple of PPT algorithms  $(\text{KG}, \text{Enc}, \text{Dec})$  such that

- (1) The key generation algorithm  $\text{KG}$  takes as input the security parameter  $\lambda$  and outputs a pair of keys  $(pk, sk)$ . For given  $pk$ , the message space  $\mathbf{M}_{pk}$  and the randomness space  $\mathbf{R}_{pk}$  are uniquely determined.
- (2) The encryption algorithm  $\text{Enc}$  takes as input a public key  $pk$  and a message  $m \in \mathbf{M}_{pk}$ , and outputs a ciphertext  $c \in \mathbf{C}_{pk}$  where  $\mathbf{C}_{pk}$  is a finite set of ciphertexts. We write this as  $c \leftarrow \text{Enc}_{pk}(m)$ . We sometimes write  $\text{Enc}_{pk}(m)$  as  $\text{Enc}_{pk}(m, r)$  when the randomness  $r \in \mathbf{R}_{pk}$  used by  $\text{Enc}$  needs to be emphasized. .
- (3) The decryption algorithm  $\text{Dec}$  takes as input a private key  $sk$  and a ciphertext  $c$ , and outputs a message  $m$  or a special symbol  $\perp$  which means failure.

We say that a public-key cryptosystem  $\mathcal{E}$  is *correct* if, for any key-pair  $(pk, sk) \leftarrow \text{KG}(\lambda)$  and any  $m \in \mathbf{M}_{pk}$ , it is the case that:  $m \leftarrow \text{Dec}_{sk}(\text{Enc}_{pk}(m))$ .

**Definition 5** ([14]) *A public-key cryptosystem  $\mathcal{E} = (\text{KG}, \text{Enc}, \text{Dec})$  with a security parameter  $\lambda$  is called to be semantically secure (IND-CPA secure) if after the standard CPA game being played with any PPT adversary  $\mathcal{A} = (\mathcal{A}_1, \mathcal{A}_2)$ , the advantage  $\text{Adv}_{\mathcal{E}, \mathcal{A}}^{\text{cpa}}(\lambda)$ , formally defined as*

$$\left| \Pr_{b,r} \left[ (pk, sk) \leftarrow \text{KG}(\lambda), (\text{state}, m_0, m_1) \leftarrow \mathcal{A}_1(pk), \right. \right. \\ \left. \left. c = \text{Enc}_{pk}(m_b; r) : b \leftarrow \mathcal{A}_2(\text{state}, m_0, m_1, c) \right] - \frac{1}{2} \right|,$$

*is negligible in  $\lambda$  for all sufficiently large  $\lambda$ .*

In the experiment above, when we allow  $\mathcal{A}_1$  to query the decryption oracle, if the advantage  $\text{Adv}_{\mathcal{E}, \mathcal{A}}^{\text{cca2}}(\lambda)$  is negligible, we say  $\mathcal{E}$  is IND-CCA1 secure, in short, CCA1 secure.