

Secure Two-Party k -Means Clustering

Paul Bunn* Rafail Ostrovsky†

Abstract

The k -Means Clustering problem is one of the most-explored problems in data mining to date. With the advent of protocols that have proven to be successful in performing *single database* clustering, the focus has changed in recent years to the question of how to extend the single database protocols to a *multiple database* setting. To date there have been numerous attempts to create specific multiparty k -means clustering protocols that protect the privacy of each database, but according to the standard cryptographic definitions of “privacy-protection,” so far all such attempts have fallen short of providing adequate privacy.

In this paper we describe a Two-Party k -Means Clustering Protocol that guarantees privacy, and is more efficient than utilizing a general multiparty “compiler” to achieve the same task. In particular, a main contribution of our result is a way to compute efficiently multiple iterations of k -means clustering without revealing the intermediate values. To achieve this, we use novel techniques to perform two-party division and sample uniformly at random from an unknown domain size.

Our techniques are quite general and can be realized based on the existence of any semantically secure homomorphic encryption scheme. For concreteness, we describe our protocol based on Paillier Homomorphic Encryption scheme (see [23]). We will also demonstrate that our protocol is efficient in terms of communication, remaining competitive with existing protocols (such as [15]) that fail to protect privacy.

Keywords: Multiparty Computation, k -means clustering

*Department of Mathematics, University of California, Los Angeles, CA 90095. Email: paulbunn@math.ucla.edu. Research partially done while visiting IPAM, and supported in part by NSF VIGRE grant DMS-0502315 and NSF Cybertrust grant no. 0430254

†Computer Science Department and Department of Mathematics, University of California, Los Angeles, CA 90095. Email: rafail@cs.ucla.edu. Research partially done while visiting IPAM, and supported in part by IBM Faculty Award, Xerox Innovation Group Award, NSF Cybertrust grant no. 0430254, and U.C. MICRO grant.

1 Introduction

1.1 Background on k -Means Clustering

The k -means clustering problem can be described as follows: A database \mathcal{D} holds information about n different objects, each object having d attributes. The information regarding each object is viewed as a coordinate in \mathbb{R}^d , and hence the objects are interpreted as data points living in d -dimensional Euclidean space. Very informally, k -means clustering algorithms are comprised of two steps. First, k initial centers are chosen in some manner, either at random or using some other “seeding” procedure. The second step is iterative (known as the “Lloyd Step”), and is described according to the following algorithm: Partition the n data points into k clusters based on which current cluster center they are closest to. Then reset the new cluster centers to be the *center of mass* (in Euclidean space) of each cluster. This process is either iterated a fixed number of times or until the new cluster centers are sufficiently close to the previous ones (based on a pre-determined measure of “sufficiently close”). The k -means clustering method is enormously popular among practitioners as an effective way to find a geometric partitioning of data points into k clusters, from which general trends or tendencies can be observed. In particular, k -means clustering is widely used in information retrieval, machine learning, and data mining research (see e.g. [21] for further discussion about the enormous popularity of k -means clustering).

The question of finding efficient algorithms for solving the k -means clustering problem has been greatly explored and is not investigated in this paper. Rather, we wish to extend an existing algorithm (which solves the k -means problem for a *single* database) to an algorithm that works in the two-database setting (in accordance with multiparty computation literature, we refer to the databases as “parties”). In particular, if two parties each hold partial data describing the d attributes of n objects, then we would like to apply this k -means algorithm to the *aggregate* data (which lives in some virtual database) in a way that protects the privacy of each party member’s data. In this paper, we will work in the most general setting, where we assume the data is *arbitrarily partitioned* between the two databases. This means that there is no assumption on how the attributes of the data are distributed among the parties (and in particular, this subsumes the case of *vertically* and *horizontally* partitioned data).

1.2 Previous Work

The k -means clustering problem is one of the functions most studied in the more general class of data-mining problems. Data-mining problems have received much attention in recent years as advances in computer technology have allowed vast amounts of information to be stored and examined. Due to

the sheer volume of inputs that are often involved in data-mining problems, generic multiparty computation (MPC) protocols become infeasible in terms of communication cost. This has led to constructions of function-specific multiparty protocols that attempt to handle a specific functionality in an efficient manner, while still providing privacy to the parties (see e.g. [18], [1], [2]).

The problem of extending *single* database k -means clustering protocols to the multiparty setting has been explored by numerous authors, whose approaches have varied widely. The main challenge in designing such a protocol is to prevent intermediate values from being leaked during the Lloyd Step. In particular, each iteration of the Lloyd Step requires k new cluster centers to be found, a process that requires division (the new cluster centers are calculated using a weighted average, which in turn requires dividing by the number of data points in a given cluster). However, the divisors should remain unknown to the parties, as leaking intermediate cluster sizes may reveal excessive information. Additionally, many current protocols for solving the *single* database k -means clustering problem improve efficiency by choosing data points according to a weighted distribution, which will then serve as preliminary “guesses” to the cluster center (e.g. [21], [4]). Choosing data points in this manner will also likely involve division.

A subtle issue that may not be obvious at first glance is how to perform these divisions in light of current cryptographic tools. In particular, most encryption schemes describe a message space that is a *finite* group (or field or ring). This means that an algorithm that attempts to solve the multiparty k -means problem in the cryptographic setting (as opposed to the information-theoretic setting) will view the data points not as elements of Euclidean Space \mathbb{R}^d , but rather as elements in \mathbb{G}^d (for some ring \mathbb{G}) in order to share encryptions of these data points with the other party members. But this then complicates the notion of “division,” which we wish to mean “division in \mathbb{R} ” as opposed to “multiplication by the inverse.” (The latter interpretation not only fails to perform the desired task of finding an *average*, but additionally may not even exist if not all elements in the ring \mathbb{G} have a multiplicative inverse).

Previous authors attempting to solve the multiparty k -means problem have incorporated various ideas to combat this obstacle. The “data perturbation” technique (e.g. [1], [2], [19]) avoids the issue altogether by addressing the multiparty k -means problem from an information-theoretical standpoint. These algorithms attempt to protect party members’ privacy by having each member first “perturb” their data (in some regulated manner), and then the perturbed data is made public to all members. Thus, the division (and all other computations) can be performed locally by each party member (on the perturbed data), and the division problem is completely avoided. Unfortunately, all current algorithms utilizing this method do not protect the privacy of the party members in the cryptographic def-

inition of privacy protection. Indeed, these protocols provide some privacy guarantee in terms of hiding the exact values of the database entries, but do not make the more general guarantee that (with overwhelming probability) *no* information can be obtained about any party’s inputs (other than what follows from the output of the function, i.e. the final cluster centers).

Another solution to the division problem (see e.g. [24]) is to have each party member perform the division locally on their own data. The problem with this method is that it requires each party to know all intermediate cluster assignments (in order to know what they should divide by), which may leak additional information and thus not satisfy complete privacy protection. A similar problem is encountered in [14], where they describe a way to privately perform division, but their protocol relies on the fact that both parties will learn the output value of the division. This means that the overall protocol is *not* secure, since the parties receive more information (the values of the quotients corresponding to intermediate divisions) than just the output of the function. One final approach, suggested by Jagannathan and Wright [15] is to interpret division as multiplication by the inverse. However, a simple example shows that this method does not satisfy correctness, i.e. does not correctly implement a k -means algorithm. (Consider e.g. dividing 11 by 5 in \mathbb{Z}_{21} . One would expect to round this to 2, but $11 \cdot 5^{-1} = 11 \cdot 17 = 19$).

One final approach encountered in the literature (see e.g. [7], [8], [9]) protects against leaking information about specific data in a different context. In this setting, the data is not distributed among many parties, but rather exists in a single database that is maintained by a trusted third party. The goal now is to have clients send requests to this third party for k -means clustering information on the data, and to ensure that the response from the server does not reveal too much information about the data. In the model we consider in this paper, these techniques cannot be applied since there is no central database or trusted third party.

To summarize, none of the existing “privacy-preserving” k -means clustering protocols provide cryptographically- acceptable security against an “honest-but-curious” adversary. We will present a formal notion of *security* in Section 2.3. Informally, the security of a multiparty protocol is defined by comparing the real-life interaction between the parties to an “ideal” scenario where a trusted third party exists. In this ideal setting, the trusted third party receives the private inputs from each of the parties, runs a k -means clustering algorithm on the aggregate data, and returns as output the final cluster centers to each party. (Note: depending on a pre-determined arrangement between the parties, the third party may also give each party the additional information of which cluster each data point belongs to.) The goal of multiparty computation is to achieve in the “real” world (where no trusted third party is assumed to exist) the same level of data privacy protection that can be obtained in the “ideal” model.

One final obstacle in designing a perfectly secure k -means clustering protocol comes from the iterative nature of the Lloyd Step. In the ideal model, the individual parties do not learn any information regarding the number of iterations that were necessary to reach the stopping condition. In the body of this paper, our main protocol will reveal this information to the parties (it is our belief that in practice, this privacy breach is unlikely to reveal meaningful information about the other party’s database). However, we discuss more fully in Appendix A alternative methods of controlling the number of iterations *without* revealing this extra information.

1.3 Our Results

We describe in Section 3 of this paper the first protocol for two-party k -means clustering that is secure against an honest-but-curious adversary (as mentioned above, general MPC protocols could in theory be applied to k -means, but any such protocol is unfeasible to use in practice). Moreover, we demonstrate that our protocol is competitive (in terms of communication and computation costs) with other current protocols (which fail to protect privacy against an honest-but-curious adversary). Exact efficiency bounds that we achieve can be found in Section 4. We remark that our honest-but-curious solution can be augmented using standard machinery (e.g., see [11, 13] and references therein) to the malicious adversary model (with a substantial increase in communication cost).

Our protocol takes as a template the *single-database* protocol of [21], and extends it to the two-party setting. We chose the particular protocol of [21] because it has two advantages over conventional single-database protocols: Firstly, it provides a provable guarantee as to the correctness of its output (assuming moderate conditions on the data); and secondly because their protocol reduces the number of iterations necessary in the Lloyd Step. However, the techniques we use to extend the single-database protocol of [21] can be readily applied to any single-database protocol.

In order to apply current cryptographic techniques to a single-database protocol in an attempt to create a secure two-party (multi-party) protocol, we are limited by the tools available today. In particular, all semantically secure homomorphic encryption schemes have a finite message space (e.g. \mathbb{Z}_N). This means that if we want to encrypt the data points (or attributes of the data points), then we must restrict the possible data values to a finite range. Therefore, instead of viewing the data points as living in \mathbb{R}^d , we “discretize” Euclidean space and approximate it via the lattice \mathbb{Z}_N^d , for some large N . All of the results of this paper are consequently restricted to the model where the data points live in \mathbb{Z}_N^d , (both in the “real” and “ideal” setting) and any function performing k -means clustering in this model is restricted to computations in \mathbb{Z}_N . Note that restricting to this “discretized” model is completely natural; indeed due to memory constraints, calculations

performed on computers are handled in this manner. As a consequence of working in the discretized space model, we also avoid privacy issues that arise from possible rounding errors (i.e. restricting input to be in \mathbb{Z}_N^d avoids the necessity of approximating inputs in \mathbb{R} by rounding up or down).

In order to extend the single database protocol of [21] to a two-party protocol, we follow the setup and some of the ideas discussed by Jagannathan and Wright in [15]. In that paper, the authors attempt to perform secure two-party k -means clustering, but (as they remark) fall short of perfect privacy due to leakage of information (including the number of data points in each cluster) that arises from an insecure division algorithm.

To solve the multiparty division problem, we define *division* in the ring \mathbb{Z}_N in a natural way, namely as the quotient Q from the Division Algorithm in the integers: $P = QD + R$. From this definition, we demonstrate how two parties can perform multiparty division in a secure manner. Additionally, we describe how two parties can select initial data points according to a weighted distribution. To accomplish this, we introduce a new protocol, the *Random Value Protocol*, which is described in Section 2.5. We note that the *Random Value Protocol* may be of independent interest as a subprotocol for other protocols that require random, oblivious sampling.

Our results utilize many existing tools and subprotocols developed in the multiparty computation literature. As such, the security guarantee of our result relies on cryptographic assumptions concerning the difficulty of inverting certain functions. In particular, we will assume the existence of a semantically secure homomorphic encryption scheme, and for ease of discussion, we use the homomorphic encryption scheme of Paillier [23].

1.4 Overview

In the next section, we briefly introduce the cryptographic tools and methods of proving privacy that we will need to guarantee security in the malicious adversary model. We also include in Section 2.2 a complete list of the subprotocols that will be used in this paper. Because most of the subprotocols that we use are general and have been described in previous MPC papers, we provide in Section 2.2 only a list of these protocols (possible implementations are included in Appendix B for completeness). An exception to this is our new *Random Value Protocol*, for which we provide a full implementation and proof of security in Section 2.5, and a description of a two-party Division Protocol in Section 2.4. Finally, in Section 3, we introduce the single database k -means clustering protocol of [21] which we then extend to a secure two-party protocol in Section 3.3.

2 Achieving Privacy

In multiparty computation (MPC) literature, devious behavior is modeled by the existence of an *adversary* who can corrupt one or more of the parties. In this paper, we will assume that the adversary is *honest-but-curious*, which means the adversary only learns the inputs/outputs of all of the corrupted parties, but the corrupted parties must engage in the protocol appropriately. We include in section 2.3 a formal definition of what it means for a protocol to “protect privacy” in the honest-but-curious adversary model (see also e.g. [11] for definitions of security against an honest-but-curious adversary).

In order to construct a private two-party k -means clustering protocol, we will utilize numerous subprotocols which themselves preserve privacy against an honest-but-curious adversary. We then use the fact that the composition of secure protocols remains secure (as proven in [5]). The novel contributions of this paper are the *Division Protocol* and the *Random Value Protocol* described in Sections 2.4 and 2.5, which are called as subprotocols in our two-party k -means clustering protocol. All of the other subprotocols that we will use perform standard functionalities, and possible implementations of these that are secure against an honest-but-curious adversary have been described by (multiple) other authors. For such functionalities, we will utilize results of other authors (and their corresponding proofs of privacy and efficiency), citing possible references. A brief description of the subprotocols that we will use can be found below in Section 2.2. Proving privacy for our two-party k -means clustering protocol will therefore be reduced to proving privacy for our two protocols. In Section 2.3 below, we classify protocols that have a specified generic form, and prove that such protocols will be secure in the honest-but-curious adversary model. Privacy of our *Division Protocol* and *Random Value Protocol* will then follow because they have this generic form. In Section 2.1 below, we first introduce the cryptographic tools we will need to guarantee privacy. The casual reader may wish to skip the description of the cryptographic tools in Section 2.1 and read only the high-level arguments of security in the first paragraph of Section 2.3, omitting the formal definitions and proofs of privacy in the rest of that section.

2.1 Cryptographic Tools

We will utilize standard cryptographic tools to maintain privacy in our two-party k -means clustering algorithm. It will be convenient to name our two participating parties, and we adopt the standard names of “Alice” and “Bob.” We will first utilize an additively homomorphic encryption scheme, e.g. Paillier ([23]). Thus, for encryptions we assume a message space \mathbb{Z}_N , where $N = pq$ is the product of two K -bit primes and K is the security parameter. In the protocols that follow, one of the parties will be responsible

for choosing the modulus N (we use the convention that Alice plays this role), and the opposite party (Bob) will be responsible for performing the requisite computations on encrypted data. The encryption scheme is a map $E : \mathbb{Z}_N \times \mathbb{H} \rightarrow \mathbb{G}$, where \mathbb{H} represents some group from which we obtain randomness, and \mathbb{G} is some other group. For notational convenience, we will write $E(m) \in \mathbb{G}$ rather than $E(m, r)$. This encryption scheme is additively homomorphic, so that: $E(m_1, r_1) + E(m_2, r_2) = E(m_1 + m_2, r_1 + r_2)$, where each addition refers to the appropriate group operation in \mathbb{G}, \mathbb{Z}_N , or \mathbb{H} . (For Paillier, $\mathbb{G} = \mathbb{Z}_{N^2}^\times$ and thus the group operation is multiplication). Additionally, the encryption scheme allows a user to (efficiently) multiply by a constant, i.e. for $c \in \mathbb{Z}_N$, anyone can compute: $cE(m, r) = E(cm, r')$. (For Paillier, if (N, g) is the public key, then $cE(m, r) := (g^{m_r N})^c = g^{mc_r cN} = g^{mc} (r^c)^N = E(cm, r')$, where $r' = r^c$).

2.2 Privacy Protecting Protocols

We list here the generic sub-protocols that will be used by our two-party k -means clustering protocol. All of the below protocols can be readily implemented using only the *Scalar Product Protocol*, and we include possible implementations in Appendix B. The *Scalar Product Protocol* is a standard protocol that has been explored much by other authors; we will not include an implementation of this protocol in this paper, but refer the reader to a number of possible references.

- *Scalar Product Protocol (SPP)*. This protocol takes in $\mathbf{x} \in \mathbb{Z}_N^t$ and $\mathbf{y} \in \mathbb{Z}_N^t$, and returns (shares of) some pre-determined degree two function $f(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^t c_i \mathbf{x}_i \mathbf{y}_i$, for public constants c_i . (See e.g. [10], where they describe such a protocol that achieves $O(tK)$ communication complexity, K the security parameter. Other implementations can be found in [22], [25] and [27].)
- *Bigger Than N Protocol (BTNP)*. Alice and Bob each hold a value in \mathbb{Z}_N . This protocol determines if the sum of their values (considered as a sum in \mathbb{Z} , not \mathbb{Z}_N) is bigger than N .
- *To Binary Protocol (TBP)*. Alice and Bob have shares of some value $X \in \mathbb{Z}_N$. If $X = x_K \dots x_1$ is the binary expansion of X , then this protocol returns shares of x_i for each $1 \leq i \leq K$. In other words, $x_i = x_i^A + x_i^B \pmod{N}$.
- *Find Minimum of 2 Numbers Protocol (FM2NP)*. Alice and Bob share two numbers. This protocol returns shares of the *location* of the smaller number (0 or 1).
- *Find Minimum of k Numbers Protocol (FMkNP)*. An extension of the above protocol, where this time as output they receive shares of the vector $(0, \dots, 1, \dots, 0)$, where the '1' appears in the m^{th} coordinate if the m^{th} number is smallest.

- *Distance Protocol (DistP)*. Computes the distance between two (shared) data points in \mathbb{Z}_N^d . An implementation of this protocol can be found in [15], which involves running the *SPP* four times on vectors of length d . Their protocol thus has communication complexity $O(dK)$.
- *Division Protocol (DivP)*. Computes the quotient (as defined below in Section 2.4) of a shared dividend by a shared divisor.
- *Computing δ_* Protocol* and *Choosing μ_1 Protocol*. These will be discussed when they arise in Sections 2.5 and 3.3.1.

2.3 Proof of Privacy

We present first the high-level argument for how our protocols will protect each party’s data. We have one of the parties (Alice) choose the encryption key, and encrypt all of her data using this key before sending it to the other party (Bob). Thus, Alice’s privacy will be guaranteed by the semantic security assumption of the encryption scheme. Meanwhile, Bob will also encrypt his data using Alice’s key, but he will blind all of the outputs he sends to Alice with randomness of his choosing, ensuring that Alice can learn nothing about his data.

We now make these notions precise by first providing a formal definition of privacy protection in the honest-but-curious adversary model, and a formal proof of privacy for the class of protocols that attempt to protect privacy in the above described manner.

Definition 1. Suppose that protocol X has Alice compute (and output) the function $f^A(\mathbf{x}, \mathbf{y})$, and has Bob compute (and output) $f^B(\mathbf{x}, \mathbf{y})$, where (\mathbf{x}, \mathbf{y}) denotes the inputs for Alice and Bob (respectively). Let $\text{VIEW}^A(\mathbf{x}, \mathbf{y})$ (resp. $\text{VIEW}^B(\mathbf{x}, \mathbf{y})$) represent Alice’s (resp. Bob’s) view of the transcript. In other words, if $(\mathbf{x}, \mathbf{r}^A)$ (resp. $(\mathbf{y}, \mathbf{r}^B)$) denotes Alice’s (resp. Bob’s) input and randomness, then:

$$\begin{aligned} \text{VIEW}^A(\mathbf{x}, \mathbf{y}) &= (\mathbf{x}, \mathbf{r}^A, m_1, \dots, m_t), \quad \text{and} \\ \text{VIEW}^B(\mathbf{x}, \mathbf{y}) &= (\mathbf{y}, \mathbf{r}^B, m_1, \dots, m_t), \end{aligned}$$

where the $\{m_i\}$ denote the messages passed between the parties. Also let $O^A(\mathbf{x}, \mathbf{y})$ and $O^B(\mathbf{x}, \mathbf{y})$ denote Alice’s (resp. Bob’s) output. Then we say that protocol X **protects privacy** (or **is secure**) against an honest-but-curious adversary if there exist probabilistic polynomial time simulators S_1 and S_2 such that:

$$\{(S_1(\mathbf{x}, f^A(\mathbf{x}, \mathbf{y})), f^B(\mathbf{x}, \mathbf{y}))\} \stackrel{c}{\equiv} \{(\text{VIEW}^A(\mathbf{x}, \mathbf{y}), O^B(\mathbf{x}, \mathbf{y}))\} \quad (1)$$

$$\{(f^A(\mathbf{x}, \mathbf{y}), S_2(\mathbf{y}, f^B(\mathbf{x}, \mathbf{y})))\} \stackrel{c}{\equiv} \{(O^A(\mathbf{x}, \mathbf{y}), \text{VIEW}^B(\mathbf{x}, \mathbf{y}))\}, \quad (2)$$

where $\stackrel{c}{\equiv}$ denotes computational indistinguishability.

With the above definition of privacy protection, we now prove the key lemma that will allow us to argue that our two-party k -means clustering protocol is secure against an honest-but-curious adversary.

Lemma 1. *Suppose that Alice has run the key generation algorithm for a semantically secure homomorphic public-key encryption scheme, and has given her public-key to Bob. Further suppose that Alice and Bob run Protocol X , for which all messages passed from Alice to Bob are encrypted using this scheme, and all messages passed from Bob to Alice are uniformly distributed (in the range of the ciphertext) and are independent of Bob's inputs. Then Protocol X is secure in the honest-but-curious adversary model.*

Proof. We prove the privacy protecting nature of Protocol X in two separate cases, depending on which party the adversary corrupts. To prove privacy, we show that for all PPT Adversaries, the view of the adversary based on Alice and Bob's interaction is indistinguishable to the adversary's view when the corrupted party interacts instead with a simulator. In other words, we show that there exist simulators S_1 and S_2 that satisfy conditions (1) and (2).

Case 1: Bob is Corrupted by Adversary. We simulate Alice's messages sent to Bob. For each encryption that Alice is supposed to send to Bob, we let the simulator S_2 pick a random element from \mathbb{Z}_N , and send an encryption of this. Any adversary who can distinguish between interaction with Alice versus interaction with S_2 can be used to break the security assumptions of E . Thus, no such PPT adversary exists, which means (2) holds.

Case 2: Alice is Corrupted by Adversary. We simulate Bob's messages sent to Alice. To do this, every time Bob is to send an encryption to Alice, the simulator picks a random element of \mathbb{Z}_N and returns an encryption of this. Again, equation (1) holds due to the fact that Alice cannot distinguish the simulator's encryption of a random number from Bob's encryption of the correct computation that has been shifted by randomness of Bob's choice. ■

2.4 Two-Party Division

As mentioned in Section 1.2, performing two-party division has been an obstacle to obtaining a secure two-party k -means clustering protocol. In this section and the next, we discuss our methods for overcoming this obstacle. In particular, we make precise what we mean by division in the ring \mathbb{Z}_N , and show that this definition not only matches our intuition as to what division should be, but also allows us to perform division in a secure way. Then in the following section, we discuss how two parties can choose a value $R \in \mathbb{Z}_Q$ uniformly at random, where $Q \in \mathbb{Z}_N$ is not known by either party, but is shared between them.

Let $P, D \in \mathbb{Z}_N$. Then viewing P and D as integers, we may apply the Division Algorithm to find unique integers $Q < N$ and $0 \leq R < D$ such that

$P = QD + R$. Viewing $Q \in \mathbb{Z}_N$, we then define *division* (of P by D) to be the quotient Q . Note that this definition is the natural restriction of division in \mathbb{R} to the integers, in that Q represents the actual quotient in \mathbb{R} that has been rounded *down* to the nearest integer. Thus this definition coincides much more closely to real division (e.g. for purposes of finding averages) than other alternatives, such as defining division to be multiplication by the inverse.

In defining what it means for a protocol to be secure (see Section 2.3), one compares the information that could be obtained in an *ideal* model (where a trusted third party exists) versus what could be obtained in the *real world* (where no such third party exists, and the proposed protocol is employed). In terms of defining the function that is to be evaluated (which performs the k -means clustering), we force the definition of division to match the above definition. In other words, when the functions $f^A(\mathbf{x}, \mathbf{y})$ and $f^B(\mathbf{x}, \mathbf{y})$ (see notation of Section 2.3) call for division to be performed, these divisions are defined to mean division in the ring \mathbb{Z}_N as defined above. This way, when our protocol is run and division is performed in this way, it matches the computations that the functions f^A and f^B are performing.

With these definitions in place, it remains to implement a secure division subprotocol that computes Q and returns shares to Alice and Bob. We describe below a possible implementation, which has been reduced to the *Scalar Product Protocol* combined with the *Find Minimum of 2 Numbers Protocol*, and consequently its security follows from the security of those subprotocols.

2.4.1 Possible Implementation of the Division Protocol

Input. Alice and Bob share $P = P^A + P^B$ and $D = D^A + D^B$ and have commitments to the other party's shares.

Output. If $P = QD + R$ for $0 \leq R < D$ is the unique expression guaranteed by the Division Algorithm, then this protocol outputs shares of Q , and commitments to these shares.

Cost. This protocol adds $O(K\xi) + O(K^3)$, where $O(\xi)$ is the communication cost of a secure *Find Minimum of K Numbers Protocol*. Note that each time the *FMKNP* is called below, the numbers are in decreasing order. As noted in Appendix B, in this case we have $O(\xi) \leq O(K^2 \log K)$.

1. Define the vector $\mathbf{D} = \mathbf{D}^A + \mathbf{D}^B$, where $\mathbf{D}^A := \{D^A, 2D^A, 2^2D^A, \dots, 2^K D^A\}$ is computed by Alice and \mathbf{D}^B is analogously computed by Bob. Note that each product involves only a single multiplication, namely by doubling the previous product.
2. Alice sets $P_0^A := P^A$ and Bob sets $P_0^B := P^B$. They then run the *Find Minimum of k Numbers Protocol (FMkNP)* on the set $(P_0, P_0 - D, P_0 - 2D, \dots, P_0 - 2^K D)$. This returns shares of $\delta_1 = \delta_1^A + \delta_1^B \in \mathbb{Z}_2^{K+1}$, which is the characteristic vector representing $a_1 \in [0..K]$, where a_1 is the

largest value such that $2^{a_1}D \leq P_0$. Define $P_1 = P_0 - 2^{a_1}D$, and notice that Alice and Bob can share $P_1 = P^A + P^B$, since $P_1 = P_0 - \delta_1 \cdot D$ (they run *SPP* to obtain shares of P_1).

3. Alice and Bob repeat the above step for $2 \leq i \leq K$. Namely, on the i^{th} iteration they run the *FMkNP* on $(P_{i-1}, P_{i-1}-D, P_{i-1}-2D, \dots, P_{i-1}-2^K D)$. This outputs δ_i , representing the characteristic vector of a_i , where a_i is the largest value in $[0..K]$ such that $2^{a_i}D < P_{i-1}$. They then obtain shares of P_i^A by running the *SPP* as in the above step.
4. Notice that $Q = (\delta_1 + \delta_2 + \dots + \delta_i) \cdot (1, 2, \dots, 2^K)$, so Alice and Bob can run the *SPP* on the appropriate function, which will output shares of Q .

2.5 The Random Value Protocol (RVP)

We describe here how two parties (Alice and Bob) can choose a value $R \in \mathbb{Z}_Q$ uniformly at random, where $Q \in \mathbb{Z}_N$ is not known by either party, but is shared between them. Before we describe the protocol, we provide motivation for why the problem is interesting. After all, with a division protocol in hand, one could simply have Alice and Bob choose an arbitrary $R' \in \mathbb{Z}_N$ (which is trivial to accomplish), and then use the division protocol to find its modulus in \mathbb{Z}_Q , and set this to be R . The problem with this approach is that if the modulus of Q in \mathbb{Z}_N is $\bar{Q} \in [0..N-1]$, then R will NOT be distributed uniformly in $[0..Q-1]$, as R will be slightly more likely to lie in $[0..\bar{Q}]$ than in $[\bar{Q}+1..Q-1]$. Since the functions f^A and f^B will be drawing R *uniformly* from \mathbb{Z}_Q , having our protocol draw R in the above way (which as noted is *not* uniformly distributed if $Q \nmid N$) will make it impossible to find simulators as in (1) and (2). We therefore need to find a way to sample *uniformly* from \mathbb{Z}_Q *without* revealing any information about Q to either party.

Recall that N is a K -bit integer, so let $Q = q_K \dots q_1$ denote the binary expansion of Q . We would like for Alice and Bob to not have any knowledge about the random value R they pick, a notion made more precise in the following definition:

Definition 2. Let VIEW^A (respectively VIEW^B) denote Alice's (resp. Bob's) view of an execution of the RVP. We say that Alice and Bob have chosen R *obliviously* if:

$$\forall Q \in \mathbb{Z}_N, \forall \alpha \in \mathbb{Z}_Q,$$

$$\Pr[R = \alpha | \text{VIEW}^A] = \Pr[R = \alpha | \text{VIEW}^B] = \frac{1}{Q}. \quad (3)$$

Additionally, we would like Q to remain unknown to both parties throughout the execution of the protocol. That Q remains unknown to both parties will follow from the fact that the below protocol is secure (as in definition 1),

and obliviousness of R will be proved in Theorem 1 below. The protocol proceeds by first describing how Alice and Bob make $S \in \mathbb{Z}_Q$ which is chosen uniformly at random, but Alice may have partial knowledge of its value (Bob however is oblivious to the value of S). This is followed by the two parties forming $T \in \mathbb{Z}_Q$ in an analogous manner but with their roles reversed, so that it is Bob who may have partial knowledge about T , and Alice who is oblivious. From these they will set $R = S + T \pmod{Q}$.

We present first a brief high-level description of how they make $S \in \mathbb{Z}_Q$. We imagine the numbers 0 through $Q - 1$ to be partitioned into groups that each have size a power of 2, as determined by the binary expansion of Q . For example, if $Q = 37 = 100101$, then we partition $[0..36]$ into the sets of size 1, 4, and 32: $\{0\}, [1..4], [5..36]$. We then choose a value from each of these sets uniformly at random, so that if there are m sets, then we choose m random values $\{S_1, \dots, S_m\}$. Finally, we set S to be one of these m values, according to a probability that depends on the size of each set. More specifically, if the i^{th} set has size 2^j , then we set S to be S_i with probability $\frac{2^j}{Q}$.

2.5.1 Description of the Protocol

Input. Alice has $Q^A \in \mathbb{Z}_N$ and Bob has $Q^B \in \mathbb{Z}_N$, and they have commitments to the other party's share of Q .

Output. Alice and Bob share $R \in \mathbb{Z}_Q$, where R has been chosen *obliviously* (as in Definition 2) and uniformly at random. More specifically, Alice has $R^A \in \mathbb{Z}_N$ and Bob has R^B such that:

$$R = R^A + R^B \pmod{N} \in [0..Q - 1].$$

Cost. This protocol will add $O(K^2)$ to communication.

Note. This protocol first requires that Bob knows the decryption key for some homomorphic encryption scheme with security parameter K , so that Alice can perform computations on their joint inputs without being able to decrypt. The protocol then flips the roles of Alice and Bob, so it is Alice who will need to hold a decryption key, and Bob who is unable to decrypt. This situation is trivial to produce, since Bob (resp. Alice) can choose their own RSA modulus N^B (resp. N^A) of K -bits, which will be used during the appropriate half of the protocol. Initially, Q is shared with respect to Alice's encryption scheme, i.e. $Q = Q^A + Q^B \pmod{N^A}$. Therefore, before running the first half of this protocol, Alice and Bob convert their shares of Q (with respect to N^A) to shares of Q (with respect to N^B). Steps 1-6 describe the first half of the protocol, where Bob's encryption key (with respect to N^B) is used, and then Step 7 (which repeats Steps 1-6 with the roles reversed) is done using Alice's encryption key (w.r.t. N^A). For ease of notation, we will drop the superscripts on N , remembering which modulus we are working in (which flips for Step 7).

1. Alice and Bob run the *To Binary Protocol (TBP)* on Q to get shares of the bits of $Q = q_K \dots q_1$.

2. Alice and Bob can now obtain shares of $Q_i = Q \pmod{2^{i-1}}$ for each $1 \leq i \leq K$ by performing the appropriate computation on their shares of the bits of Q . For instance, Alice will set:

$$Q_i^A = \sum_{j=1}^{i-1} q_j^A 2^{j-1},$$

where Q_1^A is initialized to zero. Bob does similarly to compute Q_i^B . Notice that $Q_i = Q_i^A + Q_i^B \pmod{N}$.

3. Alice picks $U \in [0..2^K - 1]$ randomly and computes U_i for each $1 \leq i \leq K$. Alice and Bob now share $S_i = U_i + Q_i = (U_i + Q_i^A) + Q_i^B \pmod{N}$. It remains to explain how Alice will pick S_i from $\{S_1, \dots, S_K\}$ with appropriate probability.
4. This step produces a reordering of $[1..K]$ such that i appears before j with probability 2^{i-j} . Label this reordering $\{x_1, \dots, x_K\}$, where each $x_i \in [1..K]$ appears exactly once. Initialize $V = 2^K - 2$ and define $V_i := V + 1 \pmod{2^{i-1}}$ for each $1 \leq i \leq K$. Alice repeats the following for each $1 \leq l \leq K$:
 - (a) Alice chooses a random number $X_l \in [0..V]$, and sets x_l to be $m \in [1..K]$ if $X_l \in [V_m..V_{m+1} - 1]$.
 - (b) Alice updates $V = V - 2^{x_l-1}$ and re-calculates each V_i .

5. This step will choose (with correct probability) the S_* (for some index $S_* \in \{S_1, \dots, S_K\}$), for which Alice will set $S = S_* = U_* + Q_*$. Namely, it will produce shares of the characteristic vector δ_* that has a ‘1’ in the $*^{th}$ coordinate and zeroes elsewhere. Letting \mathbf{e}_i denote the characteristic vector with a ‘1’ in the i^{th} position, we use the following equation to define δ_* (we leave it to the reader to verify that δ_* will choose S_* from $\{S_1, \dots, S_K\}$ with correct probability):

$$\begin{aligned} \delta_* = & (q_{x_1})\mathbf{e}_{x_1} + (1 - q_{x_1})(q_{x_2})\mathbf{e}_{x_2} + \dots + \\ & (1 - q_{x_1})(1 - q_{x_2}) \dots (1 - q_{x_{K-1}})(q_{x_K})\mathbf{e}_{x_K}. \end{aligned}$$

For brevity, we have Alice compute δ_* by running the sub-protocol *Compute δ_* Protocol*, which can be found in Appendix B with the other sub-protocols.

6. Alice and Bob can now share $S = S_*$ by running the *SPP* on the function:

$$f(\mathbf{x}, \mathbf{y}) = \delta_* \cdot (S_1, \dots, S_K).$$

7. Alice and Bob repeat steps 1-6 with their roles reversed, so that Alice and Bob share T . Now S and T are elements of \mathbf{Z}_Q , and we would like to perform the sum $S + T \pmod{Q}$. However, Alice and Bob cannot simply add their own shares of S and T because these shares correspond to two different moduli N^A and N^B . (Recall that S was created using Bob’s encryption key, and is therefore shared between

Alice and Bob modulo N^B , while T is shared between them modulo N^A .) A little work must be done to convert the shares of S (which are w.r.t. N^B) to shares of T (now w.r.t N^A), and then compute $S + T \pmod{Q}$. We leave the details to the reader.

2.5.2 Proof of Obliviousness and Security

Notice that the only communication between Alice and Bob in the above protocol takes place in the form of the sub-protocols *TBP*, *Compute δ_* Protocol*, and *SPP*. The protocol is therefore secure if each of those sub-protocols are secure, by the composition theorem of [5]. Since we are using a secure *Scalar Product Protocol* (e.g. of [10] or [25]) and the *TBP* and *Compute δ_* Protocol* (see Appendix B) are both secure, it follows that our *Random Value Protocol* is secure against an honest-but-curious adversary. It remains to show that the output $R \in \mathbb{Z}_Q$ is *oblivious* to both parties.

Theorem 1. *The above described Random Value Protocol outputs shares of $R \in \mathbb{Z}_Q$ such that R has been chosen obliviously (as in definition 2).*

Proof. The fact that R is chosen *obliviously* follows from three simple claims:

Claim 1. *During Alice’s portion of the protocol (Steps 1-6), the distribution of choices for S is uniform in \mathbb{Z}_Q . Conversely for T during Bob’s portion of the protocol (Step 7).*

Claim 2. *If β is **any** fixed number in \mathbb{Z}_Q and X represents a random variable uniformly distributed in \mathbb{Z}_Q , then the random variable $Y := \beta + X \pmod{Q}$ is uniformly distributed in \mathbb{Z}_Q .*

Claim 3. *If a party’s view includes knowledge of β but no knowledge of X , then Y is oblivious to that party.*

We leave the proofs of these claims to the reader, but note that all three claims result from straightforward combinatorial arguments. The fact that $R = S + T \pmod{Q}$ is a random variable follows from the fact that both S and T are chosen uniformly at random in \mathbb{Z}_Q , and then letting e.g. $X = S$ and $\beta = T$ in Claim 2 above, we have by Claim 2 that R is uniformly distributed in \mathbb{Z}_Q . The fact that R is oblivious follows from Claim 3. ■

As an aside, we note that Claim 2 actually guarantees that this protocol chooses R obliviously even if one of the parties is corrupted maliciously. The *Random Value Protocol* can therefore be used as a sub-protocol in models allowing a *malicious* adversary, provided that the *TBP*, *Compute δ_* Protocol*, and *SPP* utilized by the *RVP* are all secure against a malicious adversary.

3 Two-Party k -Means Clustering Protocol

3.1 Notation and Preliminaries

Following the setup of [15], we assume that two parties, “Alice” and “Bob,” each hold (partial) data describing the d attributes of n objects (we assume Alice and Bob both know d and n). Their aggregate data comprises the (virtual) database \mathcal{D} , holding the complete information of each of the n objects. The goal is to design an efficient algorithm that allows Alice and Bob to perform k -means clustering on their aggregate data in a manner that protects their private data.

As mentioned in the Introduction, we are working in the model where our data points are viewed as living in \mathbb{Z}_N^d for some large RSA modulus N chosen by Alice. Note that if Alice and Bob desire a lattice width of W and M denotes the maximum *Euclidean* distance between points, then Alice will pick N sufficiently large to guarantee that $N \geq \frac{n^2 M^2}{W^2}$ (this inequality guarantees that the sum of all data points does not exceed N).

We allow the data points to be *arbitrarily partitioned* between Alice and Bob (see [15]). This means that there is no assumed pattern to how Alice and Bob hold attributes of different data points (in particular, this subsumes the cases of *vertically* and *horizontally* partitioned data). We only demand that between them, each of the d attributes of all n data points is known by either Alice or Bob, but not both. For a given data point $\mathbf{D}_i \in \mathcal{D}$, we denote Alice’s share of its attributes by \mathbf{D}_i^A , and Bob’s share by \mathbf{D}_i^B .

3.2 Single Database k -Means Algorithms

The single database k -means clustering algorithm that we extend to the two-party setting was introduced by [21] and is summarized below. We chose this algorithm because under appropriate conditions on the distribution of the data, the algorithm is provably correct (as opposed to most other algorithms that are used in practice which have no such provable guarantee of correctness). Additionally, the Initialization Phase (or “seeding process”) is done in an optimized manner, reducing the number of iterations required in the Lloyd Step. The algorithm is as follows (see [21] for details):

Step I: Initialization. This procedure chooses the cluster centers $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$ according to (an equivalent version of) the protocol described in [21]:

- A. *Center of Gravity.* Compute the center of gravity of the n data points and denote this by \mathbf{C} :

$$\mathbf{C} = \frac{\sum_{i=1}^n \mathbf{D}_i}{n} \tag{4}$$

- B. *Distance to Center of Gravity.* For each $1 \leq i \leq n$, compute the distance (squared) between \mathbf{C} and \mathbf{D}_i . Denote this as $\tilde{C}_i^0 = \tilde{C}_i^0 = \text{Dist}^2(\mathbf{C}, \mathbf{D}_i)$.

C. *Average Squared Distance.* Compute the average squared distance $\bar{C} := \frac{\sum_{i=1}^n \tilde{C}_i^0}{n}$.

D. *Pick First Cluster Center.* Pick $\boldsymbol{\mu}_1 = \mathbf{D}_i$ with probability:

$$\Pr[\boldsymbol{\mu}_1 = \mathbf{D}_i] = \frac{\bar{C} + \tilde{C}_i^0}{2n\bar{C}}. \quad (5)$$

E. *Iterate to Pick the Remaining Cluster Centers.* Pick $\boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_k$ as follows: Suppose $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_{j-1}$ have already been chosen (initially $j=2$), then we pick $\boldsymbol{\mu}_j$ by:

1. For each $1 \leq i \leq n$, calculate \tilde{C}_i^{j-1} , the distance (squared) between \mathbf{D}_i and $\boldsymbol{\mu}_{j-1}$.
2. For each $1 \leq i \leq n$, let \tilde{C}_i denote the minimum of $\{\tilde{C}_i^l\}_{l=0}^{j-1}$.
3. Update \bar{C} to be the average of \tilde{C}_i (over all $1 \leq i \leq n$).
4. Set $\boldsymbol{\mu}_j = \mathbf{D}_i$ with probability:

$$\Pr[\boldsymbol{\mu}_j = \mathbf{D}_i] = \frac{\tilde{C}_i}{n\bar{C}}.$$

Step II: Lloyd Step. Repeat the following until $\boldsymbol{\nu}_1, \dots, \boldsymbol{\nu}_k$ is sufficiently close to $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$:

- A. *Finding the Closest Cluster Centers.* For each data point $\mathbf{D}_i \in \mathcal{D}$, find the closest cluster center $\boldsymbol{\mu}_j \in \{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k\}$, and assign data point \mathbf{D}_i to cluster j .
- B. *Calculating the New Cluster Centers.* For each cluster j , calculate the new cluster center $\boldsymbol{\nu}_j$ by finding the average position of all data points in cluster j . Share these new centers between Alice and Bob as $\boldsymbol{\nu}_1^A, \dots, \boldsymbol{\nu}_k^A$ and $\boldsymbol{\nu}_1^B, \dots, \boldsymbol{\nu}_k^B$, respectively.
- C. *Checking the Stopping Criterion.* Compare the old cluster centers to the new ones. If they are “close enough,” then the algorithm returns the final cluster centers to Alice and Bob. Otherwise, Step II is repeated after reassigning the cluster centers.
- D. *Reassigning New Cluster Centers.* To reassign new cluster centers, set:

$$\begin{aligned} \boldsymbol{\mu}_1^A, \dots, \boldsymbol{\mu}_k^A &= \boldsymbol{\nu}_1^A, \dots, \boldsymbol{\nu}_k^A, \quad \text{and} \\ \boldsymbol{\mu}_1^B, \dots, \boldsymbol{\mu}_k^B &= \boldsymbol{\nu}_1^B, \dots, \boldsymbol{\nu}_k^B. \end{aligned}$$

3.3 Our Two-Party k -Means Clustering Protocol

We now extend the k -means algorithm of [21] to a two-party setting. Section 3.3.1 below discusses how to implement Step I of the above algorithm (the Initialization), and section 3.3.2 discusses how to implement Step II of the algorithm (the Lloyd Step). We discuss in Appendix A alternative approaches in the number of iterations allowed in the Lloyd Step, and why this question is an issue in terms of protecting privacy.

3.3.1 Step I: Initialization

We now describe how to extend Step I of the above algorithm to the two-party setting. In particular, we need to explain how to perform the computations from Step I in a secure way. As output, Alice should have shares of the cluster centers $\boldsymbol{\mu}_1^A, \dots, \boldsymbol{\mu}_k^A$, and Bob should have $\boldsymbol{\mu}_1^B, \dots, \boldsymbol{\mu}_k^B$, such that $\boldsymbol{\mu}_i^A + \boldsymbol{\mu}_i^B = \boldsymbol{\mu}_i$. Below we follow Step I of the algorithm from Section 3.3.1 and describe how to privately implement each step.

A. *Center of Gravity*. To implement Step A of our algorithm, we need Alice and Bob to compute and share:

$$\mathbf{C} = \frac{1}{n} \sum_{i=1}^n \mathbf{D}_i = \frac{1}{n} \sum_{i=1}^n \mathbf{D}_i^A + \frac{1}{n} \sum_{i=1}^n \mathbf{D}_i^B. \quad (4)$$

Note that the division by n in (4) should be performed in \mathbb{R} (as opposed to \mathbb{Z}_N), which is handled by the *Division Protocol (DivP)*.

1. For each $1 \leq j \leq d$, Alice and Bob run the *SPP* on inputs $\mathbf{x} = \{\mathbf{D}_{i,j}^A\}_{i=1}^n$ and $\mathbf{y} = \{\mathbf{D}_{i,j}^B\}_{i=1}^n$, and the function $f(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n \mathbf{D}_{i,j}^A + \sum_{i=1}^n \mathbf{D}_{i,j}^B$. As output to this call, Alice gets O_j^A and Bob gets O_j^B , where:

$$O_j^A + O_j^B = O_j := \sum_{i=1}^n \mathbf{D}_{i,j}^A + \sum_{i=1}^n \mathbf{D}_{i,j}^B.$$

2. For each $1 \leq j \leq d$, Alice and Bob run the *DivP* on inputs $X^A = O_j^A$, $X^B = O_j^B$, and $D := n$. Note that as output of the *DivP*, Alice and Bob share $\mathbf{C}_j = \frac{1}{n} \sum_{i=1}^n \mathbf{D}_{i,j} = (j^{\text{th}} \text{ coordinate of } \mathbf{C})$ as desired.

B. *Distance to Center of Gravity*.

1. For each $1 \leq i \leq n$, Alice and Bob run the *Distance Protocol (DistP)* on $(\mathbf{D}_i^A, \mathbf{C}^A)$ and $(\mathbf{D}_i^B, \mathbf{C}^B)$ to share $\tilde{C}_i^0 = \tilde{C}_i^{A,0} + \tilde{C}_i^{B,0}$.

C. *Average Squared Distance*. Define the following sums:

$$P := \sum_{i=1}^n \tilde{C}_i^{A,0} \quad \text{and} \quad P' := \sum_{i=1}^n \tilde{C}_i^{B,0}$$

1. Alice and Bob run the *SPP* on inputs $\mathbf{x} = \{\tilde{C}_i^{A,0}\}_{i=1}^n$, $\mathbf{y} = \{\tilde{C}_i^{B,0}\}_{i=1}^n$, and function $f(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n \tilde{C}_i^{A,0} + \sum_{i=1}^n \tilde{C}_i^{B,0}$. As output to this function, Alice and Bob share:

$$X^A + X^B = P + P' = \sum_{i=1}^n \tilde{C}_i^{A,0} + \sum_{i=1}^n \tilde{C}_i^{B,0}.$$

2. Alice and Bob next run the *DivP* on the inputs X^A and X^B , and $D := n$. As output, Alice and Bob will be sharing \bar{C}^0 as desired.

D. *Pick First Cluster Center*. Notice that picking a data point \mathbf{D}_i with probability $\frac{\bar{C} + \tilde{C}_i^0}{2n\bar{C}}$ is equivalent to picking a random number $R \in [0..2n\bar{C} - 1]$

and finding the first i such that $R \leq \sum_{j=1}^i \bar{C} + \tilde{C}_j^0$. We use this observation to pick data points according to weighted probabilities as follows:

1. *Picking a Random R .* In this step, Alice and Bob pick a random number in $[0..2n\bar{C} - 1]$, where $2n\bar{C} = 2n\bar{C}^A + 2n\bar{C}^B$. Alice and Bob run the *Random Value Protocol (RVP)* with $Q := 2n\bar{C} = 2n\bar{C}^A + 2n\bar{C}^B$ to generate and share a random number $R = R^A + R^B \in \mathbb{Z}_{2n\bar{C}}$.
2. Alice and Bob will next compare their random number R with the sum $\sum_{j=1}^i \bar{C} + \tilde{C}_j^0$, and find the first i such that $R \leq \sum_{j=1}^i \bar{C} + \tilde{C}_j^0$. They will then set $\boldsymbol{\mu}_1 = \mathbf{D}_i$. This essentially boils down to running the *FM2NP* n times, and looking for the first place it returns a 1. The actual implementation of this can be found in the *Choosing $\boldsymbol{\mu}_1$ Protocol* in Appendix B.

E. *Iterate to Pick the Remaining Cluster Centers.*

1. This step is done analogously to Step I.B.
2. This step is supposed to calculate the minimum of $\{\tilde{C}_i^l\}_{l=0}^{j-1}$. However, they don't have to take the minimum over all j numbers, since from the previous iteration of this step, they already have $\tilde{C}_i = \text{Min}\{\tilde{C}_i^l\}_{l=0}^{j-2}$. Thus, they really only need to take a minimum of two numbers, that is reset \tilde{C}_i to be:

$$\tilde{C}_i = \text{Min}\{\tilde{C}_i, \tilde{C}_i^{j-1}\}.$$

Therefore, Alice and Bob run the *FM2NP* on inputs $(\tilde{C}_i^A, \tilde{C}_i^{A,j-1})$ and $(\tilde{C}_i^B, \tilde{C}_i^{B,j-1})$ so that they share *the location of* (the new) \tilde{C}_i (let $L = L^A + L^B$ denote this location). They can then share the new $\tilde{C}_i = \text{Min}\{\tilde{C}_i, \tilde{C}_i^{j-1}\}$ by running the *SPP* on inputs $\mathbf{x} = (\tilde{C}_i^A, \tilde{C}_i^{A,j-1}, L^A)$ and $\mathbf{y} = (\tilde{C}_i^B, \tilde{C}_i^{B,j-1}, L^B)$ and function $f(\mathbf{x}, \mathbf{y}) = L\tilde{C}_i^{j-1} + (1 - L)\tilde{C}_i$.

3. This step is done analogously to Step I.C.
4. This step is done analogously to Step I.D.

3.3.2 Step II: Lloyd Step

In this section, we discuss how to implement the Lloyd Step while maintaining privacy protection.

A. *Finding the Closest Cluster Centers.* We repeat the following procedure for each $\mathbf{D}_i \in \mathcal{D}$:

1. *Find the Distance (squared) to Each Cluster Center.* Note that because finding the minimum of all distances is equivalent to finding the

minimum of the distances squared, we will calculate the latter. Alice and Bob run the *Distance Protocol (DistP)* k times (once for each cluster j) to generate:

$$\mathbf{X}_i^A := (\mathbf{X}_{i,1}^A, \dots, \mathbf{X}_{i,k}^A) \text{ and } \mathbf{X}_i^B := (\mathbf{X}_{i,1}^B, \dots, \mathbf{X}_{i,k}^B),$$

where

$$\mathbf{X}_{i,j}^A + \mathbf{X}_{i,j}^B = \text{DistP}(\mathbf{D}_i, \boldsymbol{\mu}_j).$$

2. Alice and Bob run the *Find Minimum of k Numbers Protocol (FMkNP)* on \mathbf{X}_i^A and \mathbf{X}_i^B to obtain a share of (a vector representation of) the location of the closest cluster center to \mathbf{D}_i :

$$\mathbf{C}_i := (0, \dots, 0, 1, 0, \dots, 0), \quad (6)$$

where the 1 appears in the j^{th} coordinate if cluster center $\boldsymbol{\mu}_j$ is closest to \mathbf{D}_i . Note that in actuality, \mathbf{C}_i is shared between Alice and Bob:

$$\mathbf{C}_i = \mathbf{C}_i^A + \mathbf{C}_i^B.$$

B. Calculating the New Cluster Centers. The following will be done for each cluster $1 \leq j \leq k$. We break the calculation into three steps: In Step 1, Alice and Bob will compute and share the sum of data points in cluster j , in Step 2 they will compute and share the total number of points in cluster j , and in Step 3 they will divide the result of Step 1 by the result of Step 2. To simplify the notation, by $E(\mathbf{D}_i)$ we will mean $(E(\mathbf{D}_{i,1}), \dots, E(\mathbf{D}_{i,d}))$.

1. *Sum of Data Points in Cluster j .* In this step, Alice and Bob compute and share the sum of all data points in cluster j . We denote this sum as:

$$\mathbf{S}_j \in \mathbb{Z}_N^d = \sum_{i=1}^n \begin{cases} \mathbf{D}_i, & \text{if } \mathbf{D}_i \in \text{cluster } j \\ 0, & \text{O.W.} \end{cases}$$

At the end of this step, Alice and Bob will share $\mathbf{S}_j = \mathbf{S}_j^A + \mathbf{S}_j^B$ (here the addition is in \mathbb{Z}_N^d). Recall from Step A above that for each data point \mathbf{D}_i , Alice and Bob have \mathbf{C}_i^A and \mathbf{C}_i^B (respectively) such that:

$$\mathbf{C}_i^A + \mathbf{C}_i^B = \mathbf{C}_i = (0, \dots, 0, 1, 0, \dots, 0),$$

where the 1 appears in the m^{th} cluster if \mathbf{D}_i is closest to cluster m . Therefore, for cluster j we would like to sum:

$$\mathbf{S}_j = \sum_{i=1}^n \mathbf{C}_{i,j} \mathbf{D}_i.$$

- (a) Alice and Bob will run the *SPP* n times, where on the i^{th} time they set $\mathbf{x} = (\mathbf{C}_{i,j}^A, \mathbf{D}_i^A, \mathbf{C}_{i,j}^A \mathbf{D}_i^A)$ and $\mathbf{y} = (\mathbf{C}_{i,j}^B, \mathbf{D}_i^B, \mathbf{C}_{i,j}^B \mathbf{D}_i^B)$ and function $f(\mathbf{x}, \mathbf{y}) = \mathbf{C}_{i,j} \mathbf{D}_i = (\mathbf{C}_{i,j}^A + \mathbf{C}_{i,j}^B)(\mathbf{D}_i^A + \mathbf{D}_i^B)$. (In order to do this, they first pre-compute $\mathbf{C}_{i,j}^B \mathbf{D}_i^B$ and $\mathbf{C}_{i,j}^B \mathbf{D}_i^A$).
- (b) Notice that \mathbf{S}_j^A is the sum of all of Alice's shares from each step of the n calls to *SPP* above (and similarly for Bob and \mathbf{S}_j^B). They can therefore add all of their individual shares of the above sums to obtain shares of \mathbf{S}_j .

2. *Number of Data Points in Cluster j .* Now Alice and Bob wish to compute and share the total number of points in cluster j , denoted by T_j . To do this, for each $1 \leq i \leq n$, Alice views $\mathbf{C}_i^A \in \mathbb{Z}_N^k$, and analogously for Bob. They then run the *SPP* on inputs $\mathbf{x} = \{\mathbf{C}_i^A\}$, $\mathbf{y} = \{\mathbf{C}_i^B\}$ and function $f(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n \mathbf{C}_{i,j}^A + \sum_{i=1}^n \mathbf{C}_{i,j}^B$. (Note: since each term in the sum is in \mathbb{Z}_N^k , they actually run the *SPP* k times, once for each coordinate.)
3. *Centroid of Data Points in Cluster j .* In this step Alice and Bob would like to divide $\mathbf{S}_j^A + \mathbf{S}_j^B$ (from Step 1) by the total number of data points T_j in cluster j to obtain the new cluster center $\boldsymbol{\nu}_j$:

$$\boldsymbol{\nu}_j = \frac{\mathbf{S}_j^A + \mathbf{S}_j^B}{T_j^A + T_j^B} \quad (7)$$

Alice and Bob run the *DivP* k times (once for each cluster j) on inputs $P = S_j^A + S_j^B$ and divisor $D = T_j^A + T_j^B$, where they know $D \in [0..n]$.

C. *Checking the Stopping Criterion.* Alice and Bob run the *DistP* k times, on the i^{th} time it outputs shares of $\|\boldsymbol{\mu}_i - \boldsymbol{\nu}_i\|^2$. They can then run the *SPP* to add their shares together and run the *FM2NP* to compare these sums with ϵ , some agreed upon predetermined value. They can open then compare their outputs from the *FM2NP* to determine if the stopping criterion has been met.

D. *Reassigning New Cluster Centers.* The final step of our algorithm, replacing the old cluster centers with the new ones, is easily accomplished:

$$\begin{aligned} \text{Alice sets: } (\boldsymbol{\mu}_1^A, \dots, \boldsymbol{\mu}_k^A) &= (\boldsymbol{\nu}_1^A, \dots, \boldsymbol{\nu}_k^A), \text{ and} \\ \text{Bob sets: } (\boldsymbol{\mu}_1^B, \dots, \boldsymbol{\mu}_k^B) &= (\boldsymbol{\nu}_1^B, \dots, \boldsymbol{\nu}_k^B). \end{aligned}$$

4 Conclusion

As mentioned in Section 2.3, the proof of security of the two-party k -means clustering protocol presented above follows from the fact that each of the subprotocols are secure. The only exception to this is in step C of the Lloyd Step, where Alice and Bob must decide if their protocol has reached the termination condition. Although Alice and Bob remain oblivious to any actual values at this stage, they will gain the information of exactly how many iterations were required in the Lloyd Step. There are various ways of defining the model to handle this potential information leak and thus maintain perfect privacy protection (see Appendix A).

Analyzing the communication between Alice and Bob in the two-party k -means clustering protocol presented in Section 3.3 demonstrates that our protocol achieves communication complexity:

$$O(kK^2) + O(mndkK) + O(mnk\xi) + O((d+m)k\zeta).$$

Recall that k is the number of clusters, K is the security parameter, n is the number of data points, d is the number of attributes of each data point, m is the number of iterations in the Lloyd Step, $O(\zeta)$ is the communication cost of performing two-party secure division (where division is defined as in Section 2.4), and $O(\xi)$ is the communication cost of (securely) finding the minimum of two numbers. In this paper, we showed that $O(\zeta) \leq O(K\xi) + O(K^3)$ and that $O(\xi) \leq O(K^2 \log K)$, which means our protocol has communication complexity bounded by $O(mndkK) + O(mnkK^2) + O((m+d)kK^3 \log K)$. Notice that the cost of performing the *single database* k -means clustering protocol of [21] is at least $O(mndk) + O(mnkK) + O(mdk\zeta)$: The first term is necessary e.g. to add together all the data points in each cluster during each iteration of the Lloyd Step, the second term is necessary to e.g. find the minimum of k numbers for each data point (when deciding which cluster the data point belongs to), and the third term is necessary for performing a division for each dimension of each cluster center. Therefore, the difference in communication complexities between our secure two-party protocol and a non-secure single database protocol is at most a factor of the security parameter K . The communication cost of our protocol matches the communication complexity of [15] while simultaneously enjoying the extra guarantee of security against an honest-but-curious adversary.

The communication complexity of our k -means protocol is bounded by the cost of performing secure division. In this paper, we defined a notion of division in the ring \mathbb{Z}_N that matches the intuition of what division “should” mean (e.g. when taking an average), and we included a possible implementation of a secure division. It is our belief that improvement of our result (in terms of communication complexity) will likely be restricted to the possibility of implementing a more efficient division protocol, which is an interesting open problem.

References

- [1] D. Agrawal and C. Aggarwal. “On the Design and Quantification of Privacy Preserving Data Mining Algorithms.” *Proc. of the 20th ACM SIGMOD-SIGACT-SIGART Symp. on Principles of Database Systems*, pp. 247-255. 2001.
- [2] R. Agrawal and R. Srikant. “Privacy-Preserving Data Mining.” *Proc. of the 2000 ACM SIGMOD Int. Conf. on Management of Data*, pp. 439-450. 2000.
- [3] D. Beaver. “Foundations of Secure Interactive Computing.” *CRYPTO '91, LNCS 576*, pp. 377-391. 1992.
- [4] P. Bradley and U. Fayyad. “Refining Initial Points for K -Means Clustering.” *Proc. of the 15th International Conference on Machine Learning*, pp. 91-99. 1998.

- [5] R. Canetti. “Security and Composition of Multiparty Cryptographic Protocols.” *Journal of Cryptology*, vol. 13 no. 1 pp. 143-202. 2000.
- [6] D. Chaum, C. Crépeau and I. Damgard. “Multiparty Unconditionally Secure Protocols.” *Proc. of the 20th Annual ACM Symp. on the Theory of Computing*, pp. 11-19. 1988.
- [7] C. Dwork, F. McSherry, K. Nissim, and A. Smith. “Calibrating Noise to Sensitivity Private Data Analysis.” *Proc. of the 3rd Theory of Cryptography Conference*, pp. 265-284. 2006.
- [8] I. Dinur and K. Nissim. “Revealing Information While Preserving Privacy.” *Proc. of the 22nd ACM SIGMOD-SIGACT-SIGART Symp. on Principles of Database Systems*, pp. 202-210. 2003.
- [9] C. Dwork and K. Nissim. “Privacy-Preserving Datamining on Vertically Partitioned Data-bases.” *CRYPTO '04, LNCS 3152*, pp. 528-544. 2004.
- [10] B. Goethals, S. Laur, H. Lipmaa and T. Mielikäinen. “On Private Scalar Product Computation for Privacy-Preserving Data Mining.” *ICISC, LNCS 3506*, pp. 104-120. 2004.
- [11] O. Goldreich. “The Foundations of Cryptography, Basic Applications.” Cambridge University Press. 2004.
- [12] O. Goldreich, S. Micali and A. Wigderson. “How to Play Any Mental Game.” *Proc. of the 19th STOC, ACM*, pp. 218-229. 1987.
- [13] Y. Ishai, E. Kushilevitz, R. Ostrovsky, and A. Sahai. “Zero-Knowledge from Secure Multiparty Computation.” *ACM Symposium on Theory of Computing*. 2007
- [14] S. Jha, L. Kruger and P. McDaniel. “Privacy Preserving Clustering.” *10th European Symp. on Research in Computer Security*, pp. 397-417. 2005.
- [15] G. Jagannathan and R. Wright. “Privacy-Preserving Distributed k -Means Clustering over Arbitrarily Partitioned Data.” *KDD '05*, pp. 593-599. 2005.
- [16] J. Katz and R. Ostrovsky. “Round-Optimal Secure Two-Party Computation.” *CRYPTO '04, LNCS 3152*, pp. 335-354. 2004.
- [17] E. Kiltz, G. Leander and J. Malone-Lee. “Secure Computation of the Mean and Related Statistics.” *TCC '05, LNCS 3378*, pp. 283-302. 2005.
- [18] Y. Lindell and B. Pinkas. “Privacy Preserving Data Mining.” *CRYPTO '00, LNCS 1880*, pp. 36-54. 2000.
- [19] S. Oliveira and O.R. Zaiane. “Privacy Preserving Clustering by Data Transformation.” *Proc. 18th Brazilian Symposium on Databases*, pp. 304-318. 2003.

- [20] M. Ben-Or, S. Goldwasser, and A. Wigderson. “Completeness Theorems for Non-Cryptographic Fault-Tolerant Distributed Computation.” *Proc. 20th Annual ACM Symp. on Theory of Computing*, pp. 1-10. 1988.
- [21] R. Ostrovsky, Y. Rabani, L. Schulman, and C. Swamy. “The Effectiveness of Lloyd-Type Methods for the k -Means Problem.” *FOCS*. 2006.
- [22] M. Naor and B. Pinkas. “Oblivious Polynomial Evaluation.” *SIAM Journal of Computing*, Vol. 35, No. 5, pp. 1254-1281. 2006.
- [23] P. Paillier. “Public Key Cryptosystems Based on Composite Degree Residuosity Classes.” *Advances in Cryptology, EUROCRYPT ’99 Proceedings, LNCS 1592*, pp. 223-238. 1999.
- [24] J. Vaidya and C. Clifton. “Privacy-Preserving k -Means Clustering over Vertically Partitioned Data.” *Proc. 9th ACM SIGDD Inter. Conf. on Knowledge Discovery and Data Mining, 206-215*. 2003.
- [25] R. Wright and Z. Yang. “Privacy-Preserving Bayesian Network Structure Computation on Distributed Heterogeneous Data.” *Proc. of the 10th ACM SIGKDD International Conf. on Knowledge Discovery and Data Mining*, pp. 713-718. 2004.
- [26] A.C.C. Yao. “How to Generate and Exchange Secrets.” *Proc. of the 27th IEEE Symp. on Foundations of Computer Science*, pp. 162-167. 1986.
- [27] H. Zhu and F. Bao. “Oblivious Scalar-Product Protocols.” *11th Australasian Conference on Information Security and Privacy, LNCS 4058*, pp. 313-323. 2006.

A Alternative Computation of the k -Means Cluster Centers

It is possible that the iterative nature of the Lloyd Step may reveal undesirable information to the two parties, in particular the number of iterations that are performed in the Lloyd Step. We suggest three different approaches to handle this privacy concern:

- *Approach 1: Reveal Number of Iterations.* If Alice and Bob agree beforehand that this minor leak of information will not compromise the privacy of their data, they can choose to run our algorithm so that this is the only privacy leak.
- *Approach 2: Set the Number of Iterations to be Proportional to n .* In general, the more data points, the more iterations are necessary to

reach the stopping condition. Based on n , one could therefore approximate the expected number of iterations that should be necessary, and fix our algorithm to perform this many iterations.

- *Approach 3: Fix the Number of Iterations to be Constant.* In [21], it is argued that if the data points enjoy certain “nice” properties, then the number of iterations is extremely small (i.e. with high probability, only 2 iterations are necessary). Thus, fixing the number of iterations to be some (small) constant will (with high probability) not result in a premature termination of the Lloyd Step (i.e. the stopping condition will likely have been reached).

Each approach has its pros and cons. Approach 1 guarantees the accuracy of the final output (as the stopping criterion has been met) in the minimal number of steps, but leaks information about how many iterations were performed. Approach 2 succeeds with high probability, but may unnecessarily affect communication complexity if the fixed number of iterations is higher than necessary. Approach 3 keeps communication minimal, but runs a higher risk of losing accuracy of the final output (i.e. if the stopping criterion hasn’t been reached after the fixed number of iterations have been completed). In the body of our paper, we assumed Approach 1, although it is trivial to modify our algorithm to implement instead Approach 2 or 3.

B Implementations of Protocols from Section 2.2

We describe here possible implementations of each of the protocols listed in Section 2.2. We provide these implementations solely for the purpose of completion, and make no claim concerning their efficiency in relation to other existing protocols that perform the same tasks. Since we need each of these protocols to be secure against an honest-but-curious adversary, we need the communication in each sub-protocol to be in the generic form of Lemma 1 or to utilize other protocols that are already known to be secure; and indeed this will be the case in each of the following.

B.1 Description of the Find Minimum of 2 Numbers Protocol

Input. As input to this protocol, Alice has $(X^A, Y^A) \in \mathbb{Z}_N^2$ and Bob has $(X^B, Y^B) \in \mathbb{Z}_N^2$

Output. As output, Alice and Bob should share:

$$L = (\text{loc. of min. of } (X, Y)) := \begin{cases} 0, & \text{if } X \leq Y \\ 1, & \text{if } X \geq Y \end{cases}$$

where if $X = Y$, then L should be 0 half of the time and 1 half of the time. (Sometimes we would instead like this protocol to output 0 *always* if $X = Y$. This modification is easily accounted for by setting r in (8) below to be 0).

Cost. Total cost of this protocol is $O(K^2)$.

Note. This protocol will be completed by performing a standard minimum comparison on the binary representations of these numbers. Let $X = c_1c_2 \dots c_M$ and $Y = d_1d_2 \dots d_M$ be the binary representations of X and Y (recall that $M = \lceil \log N \rceil$). In general, note that the following formula will return the location of the minimum of (X, Y) , where the formula returns 0 if $X < Y$, a 1 if $X > Y$, and a random $r \in \{0, 1\}$ if $X = Y$:

$$\begin{aligned} L = & (c_1 \oplus d_1)c_1 + (c_1 \oplus d_1 \oplus 1)(c_2 \oplus d_2)c_2 + \\ & (c_1 \oplus d_1 \oplus 1)(c_2 \oplus d_2 \oplus 1)(c_3 \oplus d_3)c_3 + \dots + \\ & (c_1 \oplus d_1 \oplus 1) \dots (c_{M-1} \oplus d_{M-1} \oplus 1)(c_M \oplus d_M)c_M + \\ & (c_1 \oplus d_1 \oplus 1) \dots (c_M \oplus d_M \oplus 1)r \end{aligned} \quad (8)$$

where \oplus signifies XOR, and the other operations are performed in \mathbb{Z}_N . Shares of L can then be obtained by running the *SPP* many times, utilizing the general fact that:

$$c \oplus d = c + d - 2cd, \quad (9)$$

where addition on the left hand side is in \mathbb{Z}_2 and on the right hand side is in \mathbb{Z}_N . We omit the specific details due to space consideration.

B.2 Description of the Find Minimum of k Numbers Protocol

This subprotocol is a simple extension of the above. If the communication cost of the *FM2NP* is $O(\xi)$, then this protocol will have communication complexity $O(k\xi)$. Furthermore, every time this protocol is called by our k -means clustering protocol the numbers are essentially already in sorted order. We can take advantage of this and reduce the cost of this subprotocol to $O(\xi \log k)$.

B.3 Description of the To Binary Protocol

Input. As input to this protocol, Alice and Bob share $X = X^A + X^B < N/2$.

Output. If $X = x_1x_2 \dots x_M$ is the binary representation for X , then as output Alice and Bob should share each bit $x_i = x_i^A + x_i^B \pmod{N}$.

Cost. Total cost of this protocol is $O(K^2)$.

Note. This protocol is made slightly more difficult due to the two possibilities:

$$X^A + X^B = \begin{cases} X^A + X^B, & \text{if } \gamma = 0 \\ X^A + X^B - N, & \text{if } \gamma = 1 \end{cases}$$

where

$$\gamma = \begin{cases} 0, & \text{if } X^A \text{ AND } X^B < N/2 \\ 1, & \text{if } X^A \text{ OR } X^B \geq N/2 \end{cases}$$

In particular, if $X^A := a_1a_2 \dots a_M$, $X^B := b_1b_2 \dots b_M$, $2^M - N = d_1d_2 \dots d_M$, then:

$$\begin{aligned} & a_1a_2 \dots a_M \\ & b_1b_2 \dots b_M \\ & + \gamma * (d_1d_2 \dots d_M) \\ \text{BIN}(X) = & x_1x_2 \dots x_M, \end{aligned} \quad (10)$$

where addition above is standard addition in \mathbb{Z}_{2^M} (performed base 2, with carry-over). We perform addition (base 2) in the usual way: start on the right and add the bits via XOR, keeping track of carry-over. Again we omit the details, but note that addition modulo 2 can be handled by using the *SPP* together with (9).

B.4 Description of the Bigger Than N Protocol

Input. As input to this protocol, Alice and Bob share $X = X^A + X^B$, where $X < N/2$.

Output. This protocol should output shares of 0 if $X^A + X^B \geq N$ (in \mathbb{Z}), and shares of 1 otherwise.

Cost. Total cost of this protocol is $O(K)$ in communication.

Note. Define:

$$\alpha := \begin{cases} 0, & \text{if } X^A < N/2 \\ 1, & \text{if } X^A \geq N/2 \end{cases} \quad \beta := \begin{cases} 0, & \text{if } X^B < N/2 \\ 1, & \text{if } X^B \geq N/2 \end{cases} \quad (11)$$

Then due to the hypothesis that $X < N/2$, it is immediate that if $O = O^A + O^B$ denotes the output of this protocol, then:

$$\begin{aligned} O &:= \begin{cases} 0, & \text{if } X^A + X^B \pmod{N} = X^A + X^B \\ 1, & \text{if } X^A + X^B \pmod{N} = X^A + X^B - N \end{cases} \\ &= \begin{cases} 0, & \text{if } \alpha \oplus \beta = 0 \\ 1, & \text{if } \alpha \oplus \beta = 1 \end{cases} \end{aligned} \quad (12)$$

Thus, viewing the left and right hand sides of the below equation as arithmetic in \mathbb{Z}_N , and the middle as arithmetic in \mathbb{Z}_2 , we have that:

$$O^A + O^B = \alpha \oplus \beta = \alpha + \beta - 2\alpha\beta. \quad (13)$$

1. Alice set α and Bob sets β according to (11).
2. Alice and Bob run the *SPP* according to (13) to obtain shares of O .

B.5 Choosing μ_1 Protocol

Input. Alice and Bob have run the *RVP*, which has returned to them shares of a random $R \in \mathbb{Z}_{2n\bar{C}}$. They also share \bar{C} and for each $1 \leq i \leq n$, they share \tilde{C}_i^0 .

Output. Alice and Bob share $\mu_1 = \mathbf{D}_i$, where \mathbf{D}_i has been chosen with the correct probability.

Cost. This protocol costs $O(ndK) + O(nK^2)$ in terms of communication.

1. For each $1 \leq i \leq n$, Alice and Bob run the *SPP* on inputs $\mathbf{x} = (\bar{C}^A, \tilde{C}_i^{A,0})$ and $\mathbf{y} = (\bar{C}^B, \tilde{C}_i^{B,0})$ and the function:

$$f(\mathbf{x}, \mathbf{y}) = \bar{C} + \tilde{C}_i^0 = \bar{C}^A + \tilde{C}_i^{A,0} + \bar{C}^B + \tilde{C}_i^{B,0}.$$

Let $O_i = \bar{C}^A + \tilde{C}_i^{A,0} + \bar{C}^B + \tilde{C}_i^{B,0} = \bar{C} + \tilde{C}_i^0$ denote the function output value on the i^{th} call, so that Alice and Bob share this as $O_i = O_i^A + O_i^B$. Let $\mathbf{z}^A = (O_1^A, \dots, O_n^A)$, and $\mathbf{z}^B = (O_1^B, \dots, O_n^B)$ so that $\mathbf{z} = \mathbf{z}^A + \mathbf{z}^B = (\bar{C} + \tilde{C}_1^0, \dots, \bar{C} + \tilde{C}_n^0)$.

2. Alice and Bob next run the *SPP* n times to compute and share:

$$\begin{aligned}\mathbf{Z} &:= \mathbf{Z}^A + \mathbf{Z}^B = (\mathbf{z}_1, (\mathbf{z}_1 + \mathbf{z}_2), \dots, (\mathbf{z}_1 + \dots + \mathbf{z}_n)) \\ &= ((\bar{C} + \tilde{C}_1^0), \dots, (n\bar{C} + \tilde{C}_1^0 + \dots + \tilde{C}_n^0)).\end{aligned}\quad (14)$$

Notice that on any call i to *SPP*, the function involves only 3 additions- the sum of their shares from the $i - 1$ call plus the sum of their shares of \mathbf{z}_i .

3. Alice and Bob run the *Find Minimum of 2 Numbers Protocol (FM2NP)* n different times. (Actually, they run the modified version so that in the case of equality, the protocol always returns a 0). On the i^{th} try they run it on (\mathbf{Z}_i^A, R^A) and (\mathbf{Z}_i^B, R^B) , i.e. they are comparing the i^{th} coordinate of \mathbf{Z} with R (recall that R is an input value). This generates the vectors $\mathbf{L}^A := (L_1^A, \dots, L_n^A)$ and $\mathbf{L}^B := (L_1^B, \dots, L_n^B)$, where $L_m^A + L_m^B$ are the values returned by the *FM2NP* on the m^{th} call to it. Thus,

$$L_m^A + L_m^B = \begin{cases} 0, & \text{if } R \leq \mathbf{Z}_m \\ 1, & \text{if } R \geq \mathbf{Z}_m \end{cases}$$

Note that if $\mathbf{L} := \mathbf{L}^A + \mathbf{L}^B$, then it has the form: $\mathbf{L} = (0, \dots, 0, 1, 1, \dots, 1)$, where the first 1 appears in the m^{th} coordinate if m is the first time $R \leq \mathbf{Z}_m$.

4. Alice and Bob modify \mathbf{L}^A in the following way (\mathbf{L}^B is modified similarly):

$$\mathbf{L}'^A := (\mathbf{L}_1^A, (\mathbf{L}_2^A - \mathbf{L}_1^A), (\mathbf{L}_3^A - \mathbf{L}_2^A), \dots, (\mathbf{L}_n^A - \mathbf{L}_{n-1}^A))$$

so that $\mathbf{L}' = \mathbf{L}'^A + \mathbf{L}'^B$ has the form:

$$\mathbf{L}' = (0, \dots, 0, 1, 0, \dots, 0),$$

where the 1 appears in the m^{th} coordinate if m is the first time that $R \leq \mathbf{Z}_m$.

5. Now Alice and Bob simply need to take the scalar product of \mathbf{L}' with $(\mathbf{D}_1, \dots, \mathbf{D}_n)$ to correctly select μ_1 . More specifically, they do this for each dimension. Define $\mathbf{D}_m = (\mathbf{D}_{1,m}, \dots, \mathbf{D}_{n,m})$. Then $\mathbf{D}_m \in \mathbb{Z}_N^n$ represents the m^{th} coordinates of the n data points. Let \mathbf{D}_m^A represent Alice's share of \mathbf{D}_m , and \mathbf{D}_m^B denote Bob's share. Alice and Bob would like to take d scalar products (once for each dimension), where on the m^{th} time they set:

$$(\mu_{1,m}^A, \mu_{1,m}^B) = \text{Scalar Product}(\mathbf{L}^A + \mathbf{L}^B, \mathbf{D}_m^A + \mathbf{D}_m^B).$$

To accomplish this, Alice first computes $\mathbf{L}^A \cdot \mathbf{D}_m^A$, and similarly Bob computes $\mathbf{L}^B \cdot \mathbf{D}_m^B$.

6. Alice and Bob compute the desired scalar product by running *SPP* on inputs $\mathbf{x} = (\mathbf{L}^A, \mathbf{D}_m^A, \mathbf{L}^A \cdot \mathbf{D}_m^A)$ and $\mathbf{y} = (\mathbf{L}^B, \mathbf{D}_m^B, \mathbf{L}^B \cdot \mathbf{D}_m^B)$ and function $f(\mathbf{x}, \mathbf{y}) = \mathbf{L} \cdot \mathbf{D}_m$. This yields as output shares of $\mu_{1,m} = \mu_{1,m}^A + \mu_{1,m}^B$ as desired.
7. Lastly, Alice and Bob set $\mu_1^A = (\mu_{1,1}^A, \dots, \mu_{1,d}^A)$ and $\mu_1^B = (\mu_{1,1}^B, \dots, \mu_{1,d}^B)$. Notice that $\mu_1 = \mu_1^A + \mu_1^B$ is exactly as it should be, that is, $\mu_1 = \mathbf{D}_i$ with the correct probability.

B.6 Compute δ_* Protocol

Input. Alice and Bob share $Q = Q^A + Q^B$, and if $Q = q_K \dots q_1$, then for each $1 \leq i \leq K$, they also share $q_i = q_i^A + q_i^B \pmod{N}$. Alice also has a reordering of the integers $[1..M]$, which is denoted $\{x_1, \dots, x_M\}$.

Output. The vector $\delta_* = (0, \dots, 1, \dots, 0)$, a unit vector with a '1' in the appropriate coordinate, has been chosen correctly (see *RVP* for precise definition of this),

and is shared between Alice and Bob.

Cost. This protocol costs $O(K^2)$ in terms of communication.

Note. In this protocol, the roles of Alice and Bob will be reversed, so that \hat{E} will represent a homomorphic encryption function that Bob can decrypt but Alice cannot.

1. Bob sends Alice $(\hat{E}(q_1^B), \dots, \hat{E}(q_K^B))$.
2. Alice picks K values at random $\{Z_1, \dots, Z_K\} \in \mathbb{Z}_N$ and (utilizing the homomorphic properties of \hat{E}) returns to Bob $(\hat{E}(q_{x_1}^B - Z_1), \dots, \hat{E}(q_{x_K}^B - Z_K))$. Notice that Alice has *rearranged* the order in which she returns things to Bob (reflecting her choices of the x_i from the above step), but Bob doesn't know the new order because Alice has blinded each term with randomness Z_i .
3. Bob decrypts each term, and multiplies them in the following indicated manner, returning to Alice:

$$\begin{aligned} &(\hat{E}(q_{x_1}^B - Z_1), \hat{E}((q_{x_1}^B - Z_1)(q_{x_2}^B - Z_2)), \dots, \\ &\hat{E}((q_{x_1}^B - Z_1)(q_{x_2}^B - Z_2) \dots (q_{x_K}^B - Z_K))). \end{aligned}$$

4. Recall that δ_* is defined by the equation:

$$\begin{aligned} \delta_* = &(q_{x_1})\mathbf{e}_{x_1} + (1 - q_{x_1})(q_{x_2})\mathbf{e}_{x_2} + \dots + \\ &(1 - q_{x_1})(1 - q_{x_2}) \dots (1 - q_{x_{K-1}})(q_{x_K})\mathbf{e}_{x_K}. \end{aligned} \quad (15)$$

Alice now utilizes the homomorphic properties of \hat{E} to calculate (an encryption of) δ_* .

5. Alice chooses new randomness and blinds δ_* with this, returning the result to Bob who can decrypt so that Alice and Bob now share δ_* .