WILEY | Hindawi

## Research Article

# Highly Secure Privacy-Preserving Outsourced $k$-Means Clustering under Multiple Keys in Cloud Computing

**Ying Zou,[1,2] Zhen Zhao,[3] Sha Shi,[4] Lei Wang,[1] Yunfeng Peng,[5] Yuan Ping [ID],[6] and Baocang Wang [ID][3]**

[1]Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai, China
[2]Department of Mathematics Teaching and Research, Shanghai Business School, Shanghai, China
[3]State Key Laboratory of Integrated Service Networks, Xidian University, Xi'an, China
[4]Engineering Research Center of Molecular and Neuro Imaging of Ministry of Education of China
 and School of Life Science and Technology, Xidian University, Xi'an, China
[5]PBC School of Finance, Tsinghua University, Beijing, China
[6]School of Information Engineering, Xuchang University, Xuchang, China

Correspondence should be addressed to Baocang Wang; bcwang79@aliyun.com

Data clustering is the unsupervised classification of data records into groups. As one of the steps in data analysis, it has been widely researched and applied in practical life, such as pattern recognition, image processing, information retrieval, geography, and marketing. In addition, the rapid increase of data volume in recent years poses a huge challenge for resource-constrained data owners to perform computation on their data. This leads to a trend that users authorize the cloud to perform computation on stored data, such as keyword search, equality test, and outsourced data clustering. In outsourced data clustering, the cloud classifies users' data into groups according to their similarities. Considering the sensitive information in outsourced data and multiple data owners in practical application, it is necessary to develop a privacy-preserving outsourced clustering scheme under multiple keys. Recently, Rong et al. proposed a privacy-preserving outsourced $k$-means clustering scheme under multiple keys. However, in their scheme, the assistant server (AS) is able to extract the ratio of two underlying data records, and key management server (KMS) can decrypt the ciphertexts of owners' data records, which break the privacy security. AS can even reduce all data records if it knows one of the data records. To solve the aforementioned problem, we propose a highly secure privacy-preserving outsourced $k$-means clustering scheme under multiple keys in cloud computing. In this paper, noncolluded cloud computing service (CCS) and KMS jointly perform clustering over the encrypted data records without exposing data privacy. Specifically, we use BCP encryption which has additive homomorphic property and AES encryption to double encrypt data records, where the former cryptosystem prevents CCS from obtaining any useful information from received ciphertexts and the latter one protects data records from being decrypted by KMS. We first define five protocols to realize different functions and then present our scheme based on these protocols. Finally, we give the security and performance analyses which show that our scheme is comparable with the existing schemes on functionality and security.

## 1. Introduction

Data clustering [1, 2] enables data records to be classified into groups according to their features, attributes, or similarities. This property leads to its significance in many fields related to data analysis, such as pattern recognition, image processing, information retrieval, geography, and marketing. Furthermore, with the explosive data received nowadays in the information era, it has been a challenge for our digital devices not only to storage but also to perform computation on such massive data. Cloud computing relieves this problem by providing a platform with high storage capacity and strong computing power. Users tend to outsource their data on the cloud and authorize the cloud server

computing ability on data. e cloud server therefore can replace users to perform some computation on the out-sourced data, such as keyword search [3], equality test [4], and outsourced data clustering [5]. It is worth noting that, in these applications, the cloud server will send the nal result to the data owner. is gives a security issue of data integrity which has been further researched in [6–11].

By outsourced data clustering which means the cloud classi es data into di erent groups according to their similarities, it is possible to e ciently detect abnormalities, segment images, and predict diseases. As a widely applied clustering method, $k$-means clustering [1] classi es data into $k$-clusters based on their distances from cluster centers. However, the sensitive information of data on the cloud platform cannot be protected by simply using $k$-means clustering. is calls for privacy-preserving outsourced $k$-means clustering, where data is classi ed without exposing the sensitive information of data.

e traditional privacy-preserving $k$-means clustering schemes [12–15] protect the data privacy by adding noises with the sacri ce of clustering accuracy. Subsequently, some symmetric and asymmetric constructions [16–18] have been

peb

*1.3. Organization.* The rest of this paper is organized as follows. In Section 2, we recall the definitions for $k$-means clustering, BCP encryption, and AES encryption. The system model and threat models are proposed in Section 3. In Section 4, five basic protocols are constructed, and we present our scheme in which the defined protocols are invoked thoroughly. The security proof and performance analysis are given in Section 5. Finally, we conclude this paper in Section 6.

## 2. Preliminaries

*2.1. Notations.* We summarize the notations used in this paper in Table 1.

*2.2. k-Means Clustering.* $k$-means clustering is an iterative algorithm that allocates $l$ data records into $k$ disjoint clusters, each of which has a center. Let $l$ $m$-dimensional data records be $\overrightarrow{d}_1, \overrightarrow{d}_2, \ldots, \overrightarrow{d}_l$ and $k$ clusters be $c_1, c_2, \ldots, c_k$, where $\overrightarrow{\mu}_1, \overrightarrow{\mu}_2, \ldots, \overrightarrow{\mu}_k$ are the centers of $k$ clusters separately. The data record $\overrightarrow{d}_i$ will be categorized into the cluster $c_j$ if $\overrightarrow{d}_i$ and $\overrightarrow{\mu}_j$ has the minimum Euclidean distance among that of $\overrightarrow{d}_i$ and all of cluster centers. In particular, the Euclidean distance of an $m$-dimensional data record $\overrightarrow{d}_i = (d_{i,1}, d_{i,2}, \ldots, d_{i,m})$ and a cluster center $\overrightarrow{\mu}_j = (\mu_{j,1}, \mu_{j,2}, \ldots, \mu_{j,m})$ can be expressed as

$$\text{Dist}\left(\overrightarrow{d}_i, \overrightarrow{\mu}_j\right) = \sum_{h=1}^{m}\left(d_{i,h} - \mu_{j,h}\right)^2. \tag{1}$$

The detailed process of $k$-means clustering is depicted as Algorithm 1. The algorithm takes as input $l$ $m$-dimensional data records $\overrightarrow{d}_1, \overrightarrow{d}_2, \ldots, \overrightarrow{d}_l$, a predefined number of clusters $k$, and a predefined max number of iterations $I$. $k$-cluster centers are firstly picked to compute the Euclidean distance with data

$\overrightarrow{\mu}_j$ is the average value of all data $c_j$ for $j \in \{1, 2, \ldots, k\}$. If the max number of iterations is clusters does not change any more, terminate the algorithm and output the $k$-clusters.

*2.3. BCP Encryption.* In this paper, we utilize the BCP encryption proposed by Bresson et al. [32] which has the additive homomorphic property and provides double decryption mechanisms. The BCP encryption consists of five algorithms as follows:

(i) Setup$(\lambda)$. Taking as input a security parameter $\lambda$, the $p, q$ of the form $p = 2p' + 1, q = 2q' + 1$ and it computes $N = pq$, where $p', q'$, $\mathbb{G} = QR_{N^2}$, the cyclic group of quadratic residues modulo $N^2$, and we have $\text{ord}(\mathbb{G}) = N\lambda(N)/2$ with $\lambda(N) = 2p'q'$. It chooses $g \in \mathbb{G}$, the order of which is $N\lambda(N)/2$, and we have $g^{\lambda(N)} \mod N^2 = 1 + \alpha N) \mod N^2$ for $\alpha \in \{1, 2, \ldots, N - 1\}$

### Table 1: Notations.

| Symbol | Meaning |
|---|---|
| $K$ | Number of clusters |
| $L$ | Number of total data records |
| $m$ | Dimension of data records |
| $I$ | Maximum number of iterations |
| $n$ | Number of data owners |
| $n_i$ | Number of data records of $i$-th owner |
| $\overrightarrow{d}_i$ | $i$-th data record $\overrightarrow{d}_i = (d_{i,1}, d_{i,2}, \ldots, d_{i,m})$ with $i \in [1, l]$ |
| $c_i$ | $i$-th cluster with $i \in [1, k]$ |
| $\|c_i\|$ | Number of data records in cluster $c_i$ |
| $\overrightarrow{\mu}_i$ | Cluster center of $i$-th cluster |
| $\text{Dist}(\overrightarrow{a}, \overrightarrow{b})$ | Euclidean distance between vectors $\overrightarrow{a}$ and $\overrightarrow{b}$ |
| $DO_i$ | $i$-th data owner with $i \in [1, n]$ |
| QC | Query client |
| CCS | Cloud computing service |
| KMS | Key management server |
| $\lambda$ | The security parameter |
| $p, q$ | Two primes of the form $p = 2p' + 1, q = 2q' + 1$ |
| $N, \lambda(N)$ | $N = pq, \lambda(N) = 2p'q'$ |
| $\mathbb{G}$ | The cyclic group of quadratic residue modulo $N^2$ |
| $g$ | $g \in \mathbb{G}$ and its order is $N\lambda(N)/2$ |
| msk | Master secret key in BCP encryption |
| pk, sk | Public key and secret key in BCP encryption |
| ask | Symmetric key used in AES encryption |
| $\delta_{i,j}$ | $= (M_1 + r_1)(M_2 + r_2)$ Scaled squared distance between $\overrightarrow{d}_i$ and $\overrightarrow{\mu}_j$ |
| $V_{n\times k}$ | Location matrix of $n$ records in $k$ clusters |

$$\text{pp} = (N, g), \text{msk} = (p', q'). \tag{2}$$

(ii) KeyGen$(\text{pp})$.
pp
chooses $a \in [1, \text{ord}(\mathbb{G})]$ and computes $h = g^a$ mod $N^2$. Note that $\mu$ is of maximal order with high probability. It sets the output public and secret key pair $(\text{pk}, \text{sk})$ as

$$\text{pk} = h = g^a, \text{sk} = a. \tag{3}$$

(iii) Enc$(\text{pk}, M)$
message $M$, the encryption algorithm randomly chooses $r \in \mathbb{Z}_N$ and generates the ciphertext $\text{CT} = (A, B)$ as

$$A = g^r \mod N^2, B = h^r(1 + mN) \mod N^2. \tag{4}$$

Specifically, we denote $\text{Enc}_{\text{pk}}(M)$ of message $M$ under the public key pk.

(iv) Dec$(\text{sk}, \text{CT})$
ciphertext $\text{CT} = (A, B)$, the decryption algorithm output the message as

$$M = \frac{B/A^a - 1 \mod N^2}{N}. \tag{5}$$

Specifically, we denote $\text{Dec}_{\text{sk}}(\text{CT})$ as the decryption of ciphertext CT under the secret key sk.

(v) $\mathrm{sDec}\,(\mathrm{msk}, \mathrm{CT})$. Taking as input the master secret key msk and a ciphertext $\mathrm{CT} = (A, B)$, the system decryption algorithm computes

$$a \bmod N = \frac{h^{\lambda(N)} - 1 \bmod N^2}{N} \cdot \alpha^{-1} \bmod N,$$

$$r \bmod N = \frac{A^{\lambda(N)} - 1 \bmod N^2}{N} \cdot \alpha^{-1} \bmod N. \tag{6}$$

Let $ar\mathrm{ord}\,(\mathbb{G}) = \gamma_1 + \gamma_2 N$; thus, $ar = \gamma_1 \bmod N$ is e ciently computable. Let $\pi$ be the inverse of $\lambda_N$. It generates the message as

$$M = \frac{(B/g^{\gamma_1})^{\lambda(N)} - 1 \bmod N^2}{N} \cdot \pi \bmod N. \tag{7}$$

Speci cally, we denote $\mathrm{sDec}_{\mathrm{msk}}\,(\mathrm{CT})$ as the decryption of ciphertext CT under the master secret key msk.

Speci cally, BCP encryption has additive homomorphic property, which means

$$\mathrm{Enc}\,(M_1) \cdot \mathrm{Enc}\,(M_2) = \mathrm{Enc}\,(M_1 + M_2). \tag{8}$$

is property will be utilized in the whole system.

*2.4. AES Encryption.* AES encryption is an e cient symmetric encryption system widely used in practical application, where the symmetric means encryption and decryption require the same key. We give the simpli ed de nition of AES as follows:

(i) AKeyGen. e sender and receiver consult the secret key $sk$ of the AES encryption system.

(ii) AEnc. e sender generates the ciphertext CT of message $M$ under the secret key $ask$ following the AES encryption algorithm. We denote it as

$$\mathrm{CT} = \mathrm{AEnc}_{\mathrm{ask}}\,(M). \tag{9}$$

(iii) ADec. e receiver decrypts the ciphertext CT with the secret key sk. We denote it as

$$M = \mathrm{Dec}_{\mathrm{ask}}\,(\mathrm{CT}). \tag{10}$$

# 3. Models

*3.1. System Model.* As shown in Figure 1, our scheme considers four types of entities, i.e., data owner (DO), cloud computing service (CCS), key management server (KMS), and query client (QC).

(i) DO: DO has limited computing power and therefore outsources its encrypted data to the cloud. Our system involves $n$ DOs, denoted as DO

## 4. Our Construction

Based on the scheme proposed by Rong et al. in [19], we construct a more secure clustering scheme. In our construction, we utilize BCP homomorphic encryption to protect the privacy security of data owners such that adversaries cannot extract any useful information about underlying data records of data owners, while AS can easily extract $M_1/M_2$ in [19]. Furthermore, AES encryption is also used to double-encrypt the data records to prevent KMS from directly extracting data records from ciphertexts sent from DO to CCS.

*4.1. Protocols.* We first define five underlying protocols to satisfy different requirements in the clustering process. To securely transfer the data records of DO to CCS, we define secure ciphertext transformation (SCT) protocol. Since the BCP encryption used in our scheme only has additive property, we build a secure multiplication (SM) protocol to realize the multiplicative property. Finally, aiming to classify the similar data records using the ciphertexts, we construct three protocols, namely, secure distance measurement (SDM) protocol, secure distance comparison (SDC) protocol, and secure minimum distance measurement (SMDM) protocol. These protocols will be invoked through our scheme.

*4.1.1. Secure Ciphertext Transformation Protocol.* Secure ciphertext transformation (SCT) protocol aims to transfer the ciphertext of message $M$ encrypted under public key $pk_x$ to a ciphertext of $M$ encrypted under public key $pk_y$ without revealing $M$. Suppose two entities in SCT protocol, i.e., Alice and Bob, Alice interacts with Bob following SCT protocol to convert $Enc_{pk_x}(M)$ to $Enc_{pk_y}(M)$. To prevent Bob from extracting the message $M$, a random number is used to blind the message from Bob. The detailed process is listed in Algorithm 2.

Taking as the input the public keys $pk_x$ and $pk_y$ and the ciphertext $Enc_{pk_x}(M)$, Alice randomly chooses $r$

Input: Alice: $pk_x, pk_y, Enc_{pk_x}(M)$
   Bob: $msk, pk_y$
Begin: Alice:
   (a) Pick a random number $r \in \mathbb{Z}_N$
   (b) Compute $Enc_{pk_x}(M + r) \leftarrow Enc_{pk_x}(M) \cdot Enc_{pk_y}(r)$
   (c) Send $Enc_{pk_x}(M + r)$ to Bob
   Bob:
   (a) Decrypt $(M + r) \leftarrow sDec_{msk}(Enc_{pk_x}(M + r))$
   (b) Compute $Enc_{pk_y}(M + r)$
   (c) Send $Enc_{pk_y}(M + r)$ to Alice
   Alice
   (a) Compute $Enc_{pk_y}(M) = Enc_{pk_y}(M + r) \cdot Enc_{pk_y}(-r)$
Output: $Enc_{pk_y}(M)$

ALGORITHM 2: SCT protocol.

Input: Alice: $Enc_{pk_x}(M_1), Enc_{pk_x}(M_2)$
   Bob: $sk_x$.
Begin: Alice:
   (a) Pick random numbers $r_1, r_2 \in \mathbb{Z}_N$
   (b) Compute $Enc_{pk_x}(M_1 + r_1) \leftarrow Enc_{pk_x}(M_1) \cdot Enc_{pk_x}(r_1)$
   (c) Compute $Enc_{pk_x}(M_2 + r_2) \leftarrow Enc_{pk_x}(M_2) \cdot Enc_{pk_x}(r_2)$
   (d) Send $Enc_{pk_x}(M_1 + r_1), Enc_{pk_x}(M_2 + r_2)$ to Bob
   Bob:
   (a) Decrypt $(M_1 + r_1) \leftarrow Dec_{sk_x}(Enc_{pk_x}(M_1 + r_1))$
   (b) Decrypt $(M_2 + r_2) \leftarrow Dec_{sk_x}(Enc_{pk_x}(M_1 + r_1))$
   (c) Compute $\ell = (M_1 + r_1)(M_2 + r_2)$
   (d) Compute $Enc_{pk_x}(\ell)$
   (e) Send $Enc_{pk_x}(\ell)$ to Alice
   Alice:
   (a) Compute $Enc_{pk_x}(M_1 \cdot M_2) = Enc_{pk_x}(M_1)^{N-r_2}$
   $Enc_{pk_x}(M_2)^{N-r_1} \cdot Enc_{pk_x}(r_1 r_2)^{N-1} \cdot Enc_{pk_x}(\ell)$
Output: $Enc_{pk_x}(M_1 \cdot M_2)$

ALGORITHM 3: SM protocol.

*4.1.4. Secure Distance Comparison Protocol.* Secure distance comparison (SDC) protocol is to determine the shorter distance between two output distances from SDM protocol. Taking as the input two distances, i.e., $(Enc_{pk_x}(\ell_{i,a}), |c_a|)$ and $(Enc_{pk_x}(\ell_{i,b}), |c_b|))$, Alice interacts with Bob to obtain the shorter one. As in [19], the difference between two differences can be expressed as

$$Enc_{pk_x}\left(Dist\left(\overrightarrow{d}_i, \overrightarrow{\mu}_a\right)\right) \cdot Enc_{pk_x}\left(Dist\left(\overrightarrow{d}_i, \overrightarrow{\mu}_b\right)\right)^{N-1}$$

$$= Enc_{pk_x}\left(Dist\left(\overrightarrow{d}_i, \overrightarrow{\mu}_a\right) - Dist\left(\overrightarrow{d}_i, \overrightarrow{\mu}_b\right)\right)$$

$$= Enc_{pk_x}\left(t_{i,a}/|c_a|^2 - t_{i,b}/|c_b|^2\right) \qquad (12)$$

$$= Enc_{pk_x}\left(\frac{|c_b|^2 t_{i,a} - |c_a|^2 t_{i,b}}{|c_a|^2 |c_b|^2}\right).$$

Since we only need to know whether $((|c_b|^2 t_{i,a} - |c_a|^2 t_{i,b})/(|c_a|^2 |c_b|^2)) > 0$ or not, it is equal to judge whether

$|c_b|^2 t_{i,a} - |c_a|^2 t_{i,b} > 0$ or not. This means, the comparison can be related to

$$Enc_{pk_x}\left(|c_b|^2 t_{i,a} - |c_a|^2 t_{i,b}\right). \qquad (13)$$

Let $\beta$ be the maximum size of messages. We have $M \in [-2^\beta + 1, 2^\beta - 1]$, which means $M \bmod N \in [1, 2^\beta - 1]$ if $M > 0$ and $M \bmod N \in [N - 2^\beta + 1, N - 1]$. Let $\eta$ be the threshold for sign judgement chosen from $[2^\beta - 1, N + 2^\beta - 1]$. To prevent Bob from obtaining distance-related information, Alice blinds the message with a random $r \in [1, \min\{N - \eta, (N - \phi N)/2^{\beta-1}\}]$ with $\phi \in \mathbb{Z}$ and satisfying

$$\begin{cases} (2^\beta - 1) \cdot r \bmod N < \eta \\ (N - 1) \cdot r \bmod N > \eta \\ (N + 1 - 2^\beta) \cdot r \bmod N > \eta \end{cases} . \qquad (14)$$

We illustrate the detailed realization in Algorithm 5. In the process, Bob cannot obtain $t_{i,a}, t_{i,b}$.

*4.1.5. Secure Minimum Distance Measurement Protocol.* Finally, we define the secure minimum distance measurement (SMDM) protocol as Algorithm 6 to choose the shortest one among given distances.

*4.2. Our Scheme.* At the beginning, the four entities in the system, i.e., data owners DOs, query client QC, cloud computing service CCS, and key management server KMS, setup the system by running the algorithms, Setup, KeyGen, and AKeyGen. DOs then run Enc and AEnc on their data records and upload to CCS separately. CCS decrypts the received ciphertexts using ADec. After receiving the clustering request from QC, CCS interacts with KMS to transform the ciphertexts encrypted under different public keys to ciphertexts encrypted under the same public key. Subsequently, CCS performs the clustering computation. Finally, CCS interacts with KMS to transfer the clustering result to QC. It is worth noting that the defined protocols are invoked through the process.

*4.2.1. System Setup.* As the setting in the system model (see Section 4.1), we have $n$ data owners $\{DO_i\}_{1 \leq i \leq n}$, cloud computing servers (CCS), key management server (KMS), and query client (QC). Before running the protocols, related entities in the system model generate their keys as follows:

(1) Taking as the input a security parameter $\lambda$, KMS runs the setup algorithm $Setup(\lambda)$ of the BCP homomorphic cryptosystem and generate the public parameter pp and master secret key msk, where *msk* is kept secret

(2) Each data owner $DO_i$ runs $KeyGen(pp)$ to generate its own public/secret key pair $(pk_i, sk_i)$, $1 \leq i \leq n$

(3) Each $DO_i$ consults with CCS a symmetric key $ask_i$ through Diffie–Hellman key exchange protocol or other methods for $1 \leq i \leq n$

(4) CCS runs the key generation algorithm $KeyGen(pp)$ to generate its public/secret key pair as $(pk_c, sk_c)$

(5) QC runs $KeyGen(pp)$ to generate its own public/secret key pair $(pk_q, sk_q)$

*4.2.2. Data Uploading.* Following the setting in Section 4.1, assume that each data owner $DO_i$ has a dataset $D_i$ which contains $n_i$ data records, and each record has $m$ attributes, and $DO_i$ encrypts $D_i$ with BCP cryptosystem first and then

AES encryption, $1 \leq i \leq n$. Finally, $DO_i$ sends the output to CCS.

(1) $DO_i$ then runs the encryption algorithm on each record $\overrightarrow{d^i_j} = (d^i_{j,1}, d^i_{j,2}, \ldots, d^i_{j,m})$, $j \in [1, n_i]$ and obtains the encrypted result as

$$\left\{ Enc_{pk_i}\left(\overrightarrow{d}^i_j\right) \right\}_{1 \leq j \leq n_i}. \tag{15}$$

(2) To prevent the privacy disclosure from KMS, data owners double-encrypt the output ciphertext with AES encryption. Each $DO_i$ computes

$$\left\{ aEnc_{ask_i}\left(Enc_{pk_i}\left(\overrightarrow{d}^i_j\right)\right) \right\}_{1 \leq j \leq n_i}, \tag{16}$$

and sends the output results to CCS.

(3) After receiving $\left\{ aEnc_{ask_i}(Enc_{pk_i}(\overrightarrow{d}^i_j)) \right\}_{1 \leq j \leq n_i}$ from $DO_i$, CCS runs the decryption algorithm aDEC with the consulted symmetric key $ask_i$ on each ciphertext to obtain

$$\left\{ aDEC_{ask}\left(aEnc_{ask}\left(Enc_{pk_i}\left(\overrightarrow{d}^i_j\right)\right)\right) \right\}_{1 \leq i \leq n_i}$$

Input: Alice: $(Enc_{pk_x}(\ _{i,a}), |c_a|), (Enc_{pk_x}(\ _{i,b}), |c_b|)), sk_y$
    Bob: $sk_x, pk_y$
Begin: Alice interacts with Bob to compute:
    (a) $Enc_{pk_x}(\ _{i,a})$    $SM(Enc_{pk_x}(\ _{i,a}), Enc_{pk_x}(|_b|^2))$
    (b) $Enc_{pk_x}(\ _{i,b})$    $SM(Enc_{pk_x}(\ _{i,b}), Enc_{pk_x}(|_a|^2))$
Alice:
    (a) Compute $Enc_{pk_x}(\ _{i,b}) = (Enc_{pk_x}(\ _{i,b}))^{N-1}$
    (b) Compute $Enc_{pk_x}(\ _{a,b}) = Enc_{pk_x}(\ _{i,a}) \cdot Enc_{pk_x}(\ _{i,b})$
    (c) Pick a random number $r$   $\mathbb{Z}_N$
Alice interacts with Bob to compute:
    (a) $Enc_{pk_x}(r\ _{a,b})$    $SM(Enc_{pk_x}(r), Enc_{pk_x}(\ _{a,b}))$
Alice:
    (a) Send $Enc_{pk_x}(r\ _{a,b})$ to Bob
Bob:
    (a) Decrypt $r\ _{a,b}$    $Dec(Enc_{pk_x}(r\ _{a,b}))$
    (b) If $r\ _{a,b} > \eta$, $sn$    $Enc_{pk_y}(1)$; otherwise, randomly choose $r$   $\mathbb{Z}_N$ satisfying $r$   1, $sn$    $Enc_{pk_y}(r)$
    (c) Send $sn$ to Alice
Alice:
    (a) If $Dec_{sk_y}(sn) = 1$, let $Enc_{pk_x}(\ _{i,min}) = Enc_{pk_x}(\ _{i,a})$, $|c_{min}| = |c_a|$; Otherwise, we have $Dec_{sk_y}(sn)$   1, let
$Enc_{pk_x}(\ _{i,min}) = Enc_{pk_x}(\ _{i,b})$, $|c_{i,min}| = |c_b|$
Output: $(Enc_{pk_x}(\ _{i,min}), |c_{i,min}|)$

ALGORITHM 5: SDC protocol..

Input: $Enc_{pk_x}(\vec{d}_i), Enc_{pk_x}(\vec{\mu}_1), Enc_{pk_x}(\vec{\mu}_2), \ldots, Enc_{pk_x}(\vec{\mu}_k)$
Begin: **for** $\alpha = 1$ to $k$
    (a) Run $SDM(Enc_{pk_x}(\vec{d}_i), Enc_{pk_x}(\vec{\mu}_\alpha))$ and obtain the output $(Enc_{pk_x}(\ _{i,\alpha}), |c_\alpha|)$
**end for**
Let $(Enc_{pk_x}(\ _{i,min}), |c_{i,min}|) = (Enc_{pk_x}(\ _{i,1}), |c_1|)$
**for** $\alpha = 2$ to $k$
    (a) Run $SDC((Enc_{pk_x}(\ _{i,min}), |c_{i,min}|), (Enc_{pk_x}(\ _{i,2}), |c_2|))$ and obtain the output $((Enc_{pk_x}(\ _{i,min}), |c_{i,min}|))$
    (b) Set $(Enc_{pk_x}(\ _{i,min}), |c_{i,min}|)$    $(Enc_{pk_x}(\ _{i,min}), |c_{i,min}|)$
**end for**
Output: $(Enc_{pk_x}(\ _{i,min}), |c_{i,min}|)$

ALGORITHM 6: SMDM protocol.

(3) By performing the SCT protocol on all the ciphertexts received from $\{DO_i\}_{1 \ i \ n}$, CCS finally obtains

$$\left\{ Enc_{pk_c}\left(\vec{d}_j^i\right) \right\}_{1 \ i \ n, 1 \ j \ n_i}. \tag{18}$$

Let $n = n_1 + n_2 + \cdots + n_l$, and denote these $n$ ciphertexts as

$$\left\{ Enc_{pk_c}\left(\vec{d}_i\right) \right\}_{1 \ i \ n}. \tag{19}$$

For simplicity, we denote $Enc_{pk_c}(\vec{d}_i)$ as $Enc(\vec{d}_i)$ in the following.

It is worth noting that the final ciphertexts are unknown to the KMS since they are blinded in the SCT protocol.

*4.2.4. Clustering Computation.* In this phase, CCS computes the clustering results with $k$ randomly chosen cluster centers $Enc(\vec{\mu}_1), Enc(\vec{\mu}_2), \ldots, Enc(\vec{\mu}_k)$ from $\left\{ Enc_{pk_c}(\vec{d}_i) \right\}_{1 \ i \ n}$.

Let   $Enc(\vec{\mu}_i) = Enc(\vec{s}_j) = (Enc(s_{j,1}), Enc(s_{j,2}), \ldots, Enc(s_{j,m}))$ and $|c_i| = 1$. CCS also outputs a matrix $V_{n \times k}$ which refers to the location in $k$ clusters of $n$ records, where $V_{i,j} = 1$ means $\vec{d}_i$ is allocated to $j$-th cluster. In addition, there is a maximum iteration time $\phi_{max}$. Let $\phi = 0$.

(1) For a data record $Enc(\vec{d}_i)$, CCS runs the SMDM protocol on it and $k$-cluster centers with the setting $pk_x = pk_c$. Finally, CCS obtains the output

$$\left( Enc(\ _{i,min}), |c_{i,min}| \right), \tag{20}$$

where $c_{i,min} = c_\alpha$. Let $V_{i,j} = 0$ for $j$   $\alpha$.

(2) For each data record $Enc(\vec{d}_i)$ where $\vec{d}_i$   $\vec{\mu}_j, 1 \ i \ n, 1 \ j \ k$, CCS runs step 1 and obtains

$$\left( Enc(\ _{i,min}), |c_{i,min}| \right), \quad 1 \ i \ n, \tag{21}$$

and the matrix $V_{n \times k}$.

(3) With the matrix $V_{n \times k}$ and data records $\text{Enc}(\vec{d}_i) = (\text{Enc}(d_{i,1}, d_{i,2}, \ldots, d_{i,m}))$, if $V_{i,j} = 1$, CCS updates $|c_j| = |c_j| + 1$ and $\text{Enc}(s_{j,\alpha})$ as

$$\text{Enc}(d_{i,\alpha}) \cdot \text{Enc}(s_{j,\alpha}) = \text{Enc}(d_{i,\alpha} + s_{j,\alpha}), \qquad (22)$$

for $1 \leq \alpha \leq m$. Finally, CCS obtains new $|c_j|$ and $\text{Enc}(\vec{s}_j) = (\text{Enc}(s_{j,1}), \text{Enc}(s_{j,2}), \ldots, \text{Enc}(s_{j,m}))$ for $1 \leq j \leq k$. Let $\phi = \phi + 1$.

(4) If $\phi < \phi_{\max}$ and the output matrix $V_{n \times k}$ is different from that in the last iteration, CCS starts a new iteration by running steps (1), (2), and (3). Otherwise, CCS outputs the final

$$\left( \text{Enc}(\vec{s}_j), |c_j| \right)_{1 \leq j \leq k}. \qquad (23)$$

*4.2.5. Result Retrieval*

(1) CCS interacts with KMS to run the SCT protocol on $\left\{ \text{Enc}_{pk_c}(\vec{s}_j) \right\}_{1 \leq j \leq k}$ with the setting $pk_x = pk_c, pk_y = pk_q$, $\text{Enc}_{pk_x}(M) = \text{Enc}_{pk_c}(\vec{s}_j)$. CCS obtains

$$\left\{ \text{Enc}_{pk_q}(\vec{s}_j) \right\}_{1 \leq j \leq k} \qquad (24)$$

and sends it and $V_{n \times k}$ to QC.

(2) QC decrypts the received $\text{Enc}_{pk_q}(\vec{s}_j)$ with its secret key $sk_q$ by computing

$$\left\{ \text{Dec}_{sk_q}\left( \text{Enc}_{pk_q}(\vec{s}_j) \right) = \vec{s}_j \right\}_{1 \leq j \leq k}. \qquad (25)$$

QC then computes the cluster centers as

$$\left\{ \frac{\vec{s}_j}{|c_j|} \right\}_{1 \leq j \leq k}, \qquad (26)$$

where $|c_j| = \sum_{i=1}^{n} V_{i,j}$.

# 5. Security and Performance Analysis

*5.1. Security Analysis.* As shown in the proposed scheme (see Section 4.2), our protocol is realized by invoking the BCP homomorphic cryptosystem, AES encryption, and the defined protocols. Upon that, the former two cryptosystems are semantic secure, and we give the security proof of the defined protocols as follows. We take the SM protocol's security proof under "Real-vs.-Ideal" framework as an example. Other protocols' security proofs are in a similar manner and we omit here.

**Theorem 1.** *SM protocol is secure.*

*Proof.* SM protocol relates to two semihonest parties, namely, Alice and Bob. Therefore, we consider both securities of SM protocol against semihonest attacker Alice $\mathcal{A}_A$ and semihonest attacker Bob $\mathcal{A}_B$. In the protocol, Alice takes

as the input $pk_x, \text{Enc}_{pk_x}(M_1), \text{Enc}_{pk_x}(M_2)$ and Bob takes as the input the corresponding secret key $sk_x$ of public key $pk_x$.

(i) Security against $\mathcal{A}_A$: In the SM protocol, the real-world view of the attacker $\mathbb{Z}_A$ includes the input $pk_x, \text{Enc}_{pk_x}(M_1), \text{Enc}_{pk_x}(M_2)$, random numbers $r_1, r_2$, $\text{Enc}_{pk_x}(\tau)$, and the output $\text{Enc}_{pk_x}(M_1 \cdot M_2)$, where $\tau = (M_1 + r_1)(M_2 + r_2)$. $\mathcal{A}_A$ tries to obtain useful information about the underlying messages, i.e., $M_1, M_2, (M_1 + r_1)(M_2 + r_2), M_1 \cdot M_2$ that are encrypted under $pk_x$. Because of the semantic security of the used BCP homomorphic cryptosystem, we have that $\mathcal{A}_A$ cannot extract any information of underlying messages except the bit length without $sk_x$. Therefore, we can construct a simulator $\mathcal{S}_A$ in the ideal world by using ciphertexts of random chosen messages. It will be computationally hard for $\mathcal{A}_A$ to distinguish the ideal world with real world because of the semantic security of the BCP homomorphic cryptosystem. We have

$$\text{Ideal}_{\mathcal{S}_A, \mathcal{A}_A} \overset{c}{\equiv} \text{Real}_{SM, \mathcal{A}_A}, \qquad (27)$$

where $\overset{c}{\equiv}$ means computationally indistinguishable.

(ii) Security against $\mathcal{A}_B$: In the protocol, $\mathcal{A}_B$ takes as the input the secret key $sk_x$ of $pk_x$ and $\text{Enc}_{pk_x}(M_1 + r_1), \text{Enc}_{pk_x}(M_2 + r_2)$. With $sk_x$, $\mathcal{A}_B$ can decrypt the ciphertexts and obtain the underlying messages $M_1 + r_1, M_2 + r_2$. However, since $r_1, r_2$ are randomly chosen by Alice, they are random numbers in the point of view of $\mathcal{A}_B$. We can then construct a simulator $\mathcal{S}_B$ in the ideal world by using ciphertexts of random chosen messages, and it will be computationally hard for $\mathcal{A}_B$ to distinguish the ideal world with the real world. We have

$$\text{Ideal}_{\mathcal{S}_B, \mathcal{A}_B} \overset{c}{\equiv} \text{Real}_{SM, \mathcal{A}_B}. \qquad (28)$$

is completes the proof of Theorem 1. □

Next, we prove that our protocol is secure by taking the process of data uploading as an example.

**Theorem 2.** *The data uploading process is secure.*

*Proof.* In the data uploading process, data owners (DOs) double-encrypt their data records with pk and ask using the BCP homomorphic cryptosystem and AES encryption separately. They then send the encrypted result to CCS who has ask but does not have the corresponding secret key sk of pk. Because of the semantic security of the BCP homomorphic cryptosystem, it is secure against semihonest CCS. Although KMS can extract the underlying messages of ciphertexts encrypted using the BCP homomorphic cryptosystem, it is also computationally hard for a semihonest KMS to obtain any information of data records with the semantic security of AES encryption. Furthermore, CCS and KMS are supposed not to collude in our scheme such that the data

TABLE 2:  e summary of schemes.

| Scheme | S/AS | Multiple data owners | Multiple keys | Cipher comparison | Security | Multidimensional data |
|--------|------|----------------------|---------------|-------------------|----------|-----------------------|
| [22] | AS | × | × | × | × | |
| [23] | S | × | × | | | |
| [26] | AS | | × | | | |
| [30] | AS | | | | × | |
| [34] | AS | | | | × | |
| [19] | AS | | | | × | |
| Ours | AS | | | | | |

S/AS: symmetric/asymmetric.

uploading process is secure against semihonest CCS and KMS separately.  is completes the proof of  eorem 2.

It is worth noting that the security of our construction is protected by the semantic security of the BCP homomorphic cryptosystem, AES encryption, and blinding with random numbers, which prevents the adversaries from obtaining any useful information from the received ciphertexts.

*5.2. Performance Analysis.* In our construction, we use the BCP homomorphic cryptosystem and AES encryption to encrypt data owners' data records to prevent the information disclosure to KMS. Compared with the underlying scheme [19] which utilizes Youn's homomorphic encryption scheme [33], our scheme therefore increases the computation cost between DOs and CCS along with the increased security.

In particular, each data owner additionally needs to interact with CCS to consult a symmetric key of AES encryption in the system setup phase. Except this, since BCP encryption has additive homomorphic property instead multiplication in Youn's encryption scheme [33], we give a secure multiplication protocol SM instead of secure addition SA in [19].  is leads to di erent invocations in other de ned protocols, which result in more computation cost.

With the sacri ce on the computation cost, our scheme achieves semantic security that adversaries cannot obtain any useful information about underlying data records, while AS can extract $M_1/M_2$ in SA protocol of [19]. Furthermore, in our scheme, KMS cannot extract the underlying data records of data owners, while KMS can realize this with its master secret key in [19].

Finally, we compare our scheme with the existing outsourced $k$-means clustering schemes [19, 22, 23, 26, 30, 34] in Table 2 on six aspects, i.e., whether the scheme is based on symmetric or asymmetric cryptosystem, whether it supports or achieves multiple data owners and multiple keys, ciphertext comparison, security, and multidimensional data. As shown in Table 2, our scheme achieves all the listed functionalities under the asymmetric cryptosystem.

## 6. Conclusions

 is paper proposed a highly secure privacy-preserving outsourced $k$-means clustering scheme on the encrypted datasets under multiple keys. We utilized BCP homomorphic encryption and AES encryption to double-encrypt the data records in the database to protect the security against semihonest cloud computing server and key management

server. Furthermore, we constructed  ve protocols, i.e., secure ciphertext transformation (SCT), secure multiplication (SM), secure distance measurement (SDM), secure distance comparison (SDC), and secure minimum distance measurement (SMDM), as the base of our scheme. In particular, SM protocol is built to achieve the homomorphic multiplicative property using BCP encryption. Finally, we proposed our scheme by invoking the de ned protocols thoroughly.  e given security and performance analysis shows that our scheme is comparable with the existing outsourced $k$-means clustering scheme on security and functionality.

## Data Availability

 e data used to support the  ndings of this study are included within the article.

## Conflicts of Interest

 e authors declare that they have no con icts of interest.

## Acknowledgments

## References

[1] A. K. Jain, M. N. Murty, and P. J. Flynn, "Data clustering: a review," *ACM Computing Surveys*, vol. 31, no. 3, pp. 264–323, 1999.

[2] A. K. Jain, "Data clustering: 50 years beyond k-means," *Pattern Recognition Letters*, vol. 31, no. 8, pp. 651–666, 2010.

[3] J. Li, Q. Wang, C. Wang, N. Cao, K. Ren, and W. Lou, "Fuzzy keyword search over encrypted data in cloud computing," in *Proceedings of the INFOCOM*, pp. 441–445, IEEE, San Diego, CA, USA, March 2010.

[4] V. Goyal, "Reducing trust in the PKG in identity based cryptosystems," in *Proceedings of the CRYPTO 2007, 27th*

*Annual International Cryptology Conference, Lecture Notes in Computer Science*, vol. 4622, pp. 430–447, Springer, Santa Barbara, CA, USA, August 2007.

[5] D. A. Davis, N. V. Chawla, N. Blumm, N. A. Christakis, and A. Barabási, "Predicting individual disease risk based on medical history," in *Proceeding of the 17th ACM Conference on Information and Knowledge Mining*, pp. 769–778, ACM, Napa Valley, CA, USA, October 2008.

[6] J. Li, H. Yan, and Y. Zhang, "Certi cateless public integrity checking of group shared data on cloud storage," *IEEE Transactions on Services Computing*, pp. 1–10, 2018.

[7] H. Yan, J. Li, and Y. Zhang, "Remote data checking with a designated veri er in cloud storage," *IEEE Systems Journal*, pp. 1–10, 2019.

[8] J. Li, H. Yan, and Y. Zhang, "E cient identity-based provable multi-copy data possession in multi-cloud storage," *IEEE Transactions on Cloud Computing*, p. 1, 2019.

[9] G. Wu, Y. Mu, W. Susilo, F. Guo, and F. Zhang, " reshold privacy-preserving cloud auditing with multiple uploaders," *International Journal of Information Security*, vol. 18, no. 3, pp. 321–331, 2019.

[10] G. Wu, Y. Mu, W. Susilo, F. Guo, and F. Zhang, "Privacy-preserving certi cateless cloud auditing with multiple users," *Wireless Personal Communications*, vol. 106, no. 3, pp. 1161–1182, 2019.

[11] G. Wu, Y. Mu, W. Susilo, and F. Guo, "Privacy-preserving cloud auditing with multiple uploaders," *Information Security Practice and Experience*, vol. 10060, pp. 224–237, 2016.

[12] A. Blum, C. Dwork, F. McSherry, and K. Nissim, "Practical privacy: the sulq framework," in *Proceedings of the ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, ACM*, pp. 128–138, Baltimore, MD, USA, June 2005.

[13] Q. Yu, Y. Luo, C. Chen, and X. Ding, "Outlier-eliminated k-means clustering algorithm based on di erential privacy preservation," *Applied Intelligence*, vol. 45, no. 4, pp. 1179–1191, 2016.

[14] J. Ren, J. Xiong, Z. Yao, R. Ma, and M. Lin, "Dplk-means: a novel di erential privacy k-means mechanism," in *Proceedings of the 2017 IEEE Second International Conference on Data Science in Cyberspace (DSC)*, pp. 133–139, Shenzhen, China, June 2017.

[15] T. Shang, Z. Zhao, Z. Guan, and J. Liu, "A DP canopy k-means algorithm for privacy preservation of hadoop platform," in *Proceedings of the CSS 2017, Lecture Notes in Computer Science*, vol. 10581, pp. 189–198, Springer, Xi'an, China, October 2017.

[16] P. Paillier, "Public-key cryptosystems based on composite degree residuosity classes," in *Proceedings of the EURO-CRYPT'99, Lecture Notes in Computer Science*, vol. 1592, pp. 223–238, Springer, Prague, Czech Republic, May 1999.

[17] R. L. Rivest, A. Shamir, and L. Adleman, "A method for obtaining digital signatures and public-key cryptosystems," *Communications of the ACM*, vol. 26, no. 1, pp. 96–99, 1983.

[18] M. Kantarcioglu and C. Clifton, "Privacy-preserving distributed mining of association rules on horizontally partitioned data," *IEEE Transactions on Knowledge and Data Engineering*, vol. 16, no. 9, pp. 1026–1037, 2004.

[19] H. Rong, H. Wang, J. Liu, J. Hao, and M. Xian, ""Outsourced k-means clustering over encrypted data under multiple keys in spark framework," in *Proceedings of the SecureComm 2017, Lecture Notes of the Institute for Computer Sciences, Social Informatics and Telecommunications Engineering*, vol. 238,

[20] Q. Zhang, L. T. Yang, Z. Chen, and P. Li, "Pphopcm: privacy-preserving high-order possibilistic c-means algorithm for big data clustering with cloud computing," *IEEE Transactions on Big Data*, 2017.

[21] Z. Brakerski, C. Gentry, and V. Vaikuntanathan, "(Leveled) fully homomorphic encryption without bootstrapping," in *Innovations in Theoretical Computer Science*, pp. 309–325, ACM, Cambridge, MA, USA, 2012.

[22] N. Almutairi, F. Coenen, and K. Dures, "K-means clustering using homomorphic encryption and an updatable distance matrix: secure third party data clustering with limited data owner interaction," in *Proceedings of the DaWa K 2017, Lecture Notes in Computer Science*, vol. 10440, pp. 274–285, Springer, Lyon, France, August 2017.

[23] J. Yuan and Y. Tian, "Practical privacy-preserving mapreduce based k-means clustering over large-scale dataset," *IEEE Transactions on Cloud Computing*, vol. 7, no. 2, pp. 568–579, 2019.

[24] O. Regev, "On Lattices, learning with errors, random linear codes, and cryptography," in *ACM Symposium on Theory of Computing*, pp. 84–93, ACM, Baltimore, MD, USA, 2005.

[25] K.-P. Lin, "Privacy-preserving kernel k-means clustering outsourcing with random transformation," *Knowledge and Information Systems*, vol. 49, no. 3, pp. 885–908, 2016.

[26] F. Rao, B. K. Samanthula, E. Bertino, X. Yi, and D. Liu, "Privacy-preserving and outsourced multi-user k-means clustering," in *Proceedings of the CIC 2015*, pp. 80–89, IEEE Computer Society, Hangzhou, China, October 2015.

[27] Y. Liu, Y. Luo, Y. Zhu, Y. Liu, and X. Li, "Secure multi-label data classi cation in cloud by additionally homomorphic encryption," *Information Sciences*, vol. 468, pp. 89–102, 2018.

[28] Z. Gheid and Y. Challal, "E cient and privacy-preserving k-means clustering for big data mining," in *Proceedings of the 2016 IEEE Trustcom/BigDataSE/ISPA*, pp. 791–798, IEEE, Tianjin, China, August 2016.

[29] C. Clifton, M. Kantarcioglu, J. Vaidya, X. Lin, and M. Y. Zhu, "Tools for privacy preserving distributed data mining," *ACM Sigkdd Explorations Newsletter*, vol. 4, no. 2, pp. 28–34, 2002.

[30] A. Peter, E. Tews, and S. Katzenbeisser, "E ciently outsourcing multiparty computation under multiple keys," *IEEE Transactions on Information Forensics and Security*, vol. 8, no. 12, pp. 2046–2058, 2013.

[31] Y. Li, Z. L. Jiang, X. Wang, S. Yiu, and J. Fang, "Outsourced privacy-preserving random decision tree algorithm under multiple parties for sensor-cloud integration," in *Proceedings of the ISPEC 2017, Lecture Notes in Computer Science*, vol. 10701, pp. 525–538, Springer, Melbourne, Australia, December 2017.

[32] E. Bresson, D. Catalano, and D. Pointcheval, "A simple public-key cryptosystem with a double trapdoor decryption mechanism and its applications," in *Proceedings of the ASIACRYPT 2003, Lecture Notes in Computer Science*, vol. 2894, pp. 37–54, Springer, Taipei, Taiwan, November 2003.

[33] T. Youn, Y. Park, C. H. Kim, and J. Lim, "An e cient public key cryptosystem with a privacy enhanced double decryption mechanism," in *Proceedings of the SAC 2005*, vol. 3897, pp. 144–158, Springer, Kingston, ON, Canada, August 2005, Lecture Notes in Computer Science.

[34] X. Liu, R. H. Deng, K.-K. R. Choo, and J. Weng, "An e cient privacy-preserving outsourced calculation toolkit with multiple keys," *IEEE Transactions on Information Forensics and Security*, vol. 11, no. 11, pp. 2401–2414, 2016.