

# Anonymous Biometric Access Control

Shuiming Ye, Ying Luo, Jian Zhao, and Sen-ching S. Cheung

{shuiming.ye, ying.luo, jian.zhao, sen-ching.cheung}@uky.edu

## Abstract

Biometric access control systems are used to control the access of resources by users based on their physical characteristics such as fingerprints and iris patterns. Biometrics patterns enable the system operator to associate the true identity of an individual user with a complete log of his/her activities after gaining access to the system. Such information may pose a significant risk to individuals' right of privacy. In this paper, we propose the concept of Anonymous Biometric Access Control (ABAC) systems for protecting user anonymity. Our proposed systems allow the server to verify the membership status of a user without knowing his/her true identity. There are two main technical contributions of our work. First, we develop a novel privacy-preserving biometric matching scheme based on homomorphic encryption to authenticate the probe in an anonymous fashion. Second, we propose a new framework called  $k$ -Anonymous Quantization (kAQ) that significantly reduces the computational complexity of the biometric matching in encrypted domain. kAQ is conceptually similar to the  $k$ -anonymity model but is used to limit the knowledge of the server about the probe to at least  $k$  maximally dissimilar items in the biometric database. Using kAQ as a preprocessing step, the search space is reduced from the entire database to  $k$  items while privacy is optimally preserved as kAQ maximizes the dissimilarity among these  $k$  items. To demonstrate our proposed framework, experiments have been conducted using a large collection of iris biometrics whose similarity is measured via a modified hamming distance. These results not only demonstrate the validity of our framework in trading-off privacy with efficiency but also illustrate the practicality in implementing secure similarity search in realistic applications.

## I. INTRODUCTION

In the last thirty years, advances in computing technologies have brought dramatic improvements in collecting, storing, and sharing personal information among government agencies and private sectors. At the same time, new forms of privacy invasion begin to enter the public consciousness. From sale of personal information to identity theft, from credit card fraud to

Contacting Author: Sen-ching S. Cheung. All authors are with University of Kentucky. A companion software to this paper is available for download at <http://www.vis.uky.edu/mialab>.

You-Tube surrendering user data [1], the number of ways that our privacy can be violated increases rapidly.

One important area of growing concern is the protection of sensitive information in various access control systems. Access control in a distributed client-server system can generally be implemented by requesting digital credentials of the user wanting to access the system. Credentials are composed of attributes that contain identifiable information about a given user. Such information can be very sensitive and uncontrolled disclosure of such attributes can result in many forms of privacy breaches. It is unsurprising that privacy protection has been a central concern in widespread deployment of access control systems, especially in many of the e-commerce applications [2].

Among the different types of access control systems, Biometric Access Control (BAC) systems pose the most direct threat to privacy. BAC systems control allocation of resources based on highly-discriminative physical characteristics of the user such as fingerprints, iris images, voice patterns or even DNA sequences. As a biometric signal is based on “who you are” rather than “what you have”, BAC systems excel in authenticating a user’s identity. While the use of biometrics enhances system security and alleviates users from carrying identity cards or remembering passwords, it creates a conundrum for privacy advocates as the knowledge of the identity makes it much harder to keep users anonymous. A curious system operator or a parasitic hacker can infer the identity of a user based on his/her biometric probe. Furthermore, as biometrics is immutable from systems to systems, it can be used by attackers to cross-correlate disparate databases and cause damages far beyond the coverage of any protection schemes for individual database systems.

A moment of thought reveals that many access control systems do not need the true identity of the user but simply require a confirmation that the user is a legitimate member. For example, an online movie vendor may have a category of “VIP” members who pay a flat monthly membership fee and can enjoy an unlimited number of movies download. While it is important to verify the VIP status of a candidate user, it is unnecessary to precisely identify who the user is. In fact, it will be appealing to customers if the vendor can provide a guarantee that it can never track their movie selections. Entry control of a large office building that hosts many companies can also benefit from such an anonymous access control system. While it is essential to restrict entry only to authorized personnel, individual companies may be reluctant to turn over sensitive

identity information to the building management. Thus a system that can validate the tenant status of a person entering the building without knowing the true identity will be valuable. Another example is a community electronic message board. Only the members of the community can sign in to the system. Once their member status are verified, they can anonymously post messages and complaints to the entire community. All the aforementioned examples can benefit from an access control system that can verify the membership status using biometric signals while keeping the identity anonymous.

In this paper, we introduce Anonymous Biometric Access Control (ABAC) to provide anonymity and access control in such a way that the system server (*Bob*) can authenticate the membership status of a user (*Alice*) but cannot differentiate Alice from any other authorized users in his database. Our scheme differs from other work in privacy protection of biometric systems which focus primarily on the security of the biometric data from improper access. Our goal is to guarantee user's anonymity while providing the safeguard of the system resources similar to other access control systems.

In this paper, we consider two technical challenges in developing an ABAC system. First, to cope with the variability of the input probe, any biometric access system needs to perform a signal matching process between the probe and all the records in the database. The challenge here lies in making the process secure so that Bob can confirm the membership status of Alice without knowing any additional information about Alice's probe. We cast this process as a secure multiparty computation problem and develop a novel protocol based on homomorphic encryption. Such a procedure prevents Bob from extracting any knowledge about Alice's probe and its similarity distances with any records in Bob's database. On the other hand, Bob can compare the distances to a similarity threshold in the encrypted domain and the comparison results are aggregated into two secret numbers shared between Bob and Alice. The secret share held by Bob prevents Alice from cheating and Alice's membership status can be verified by Bob without knowing her identity.

Second, we consider the complexity challenge posed by scaling the matching process in encrypted domain to large databases. The high complexity of cryptographic primitives is often cited as the major obstacle of their widespread deployment in realistic systems. This is particularly true for biometric applications that require matching a large number of high-dimensional feature vectors in real time. In this paper, we propose a novel framework to provide a controllable

trade-off between privacy and complexity. We call the framework  $k$ -anonymous ABAC system ( $k$ -ABAC) which keeps Alice anonymous from  $k$ , rather than the entire database of, authorized members in the database. This is similar to the well-known  $k$ -anonymity model [3] in that  $k$  is a controllable parameter of anonymity. However, the two approaches are fundamentally different – the  $k$ -anonymity model is a data disclosure protocol where Bob anonymizes the database for public release by grouping all the data into  $k$ -member clusters. In a  $k$ -ABAC system, the goal is to *prevent Bob from obtaining information about the similarity relationship between his data and the query probe from Alice*. In order to minimize the knowledge revealed by any  $k$ -member cluster, we propose a novel grouping scheme called  $k$ -Anonymous Quantization (kAQ) that optimizes the *dissimilarity* among members in the same group. kAQ forbids similar patterns to be in the same group which might be a result of multiple registrations of the same person or from family members with similar biometric features. The kAQ process is carried out mostly in plaintext and is computationally efficient. Using kAQ as a pre-processing step, the subsequent encrypted-domain matching can be efficiently realized within the real-time constraint.

The rest of the paper is organized as follows: after reviewing related work in Section II, we provide the necessary background in the security models for anonymous biometric matching, homomorphic encryption and dimension reduction in Section III. We first provide an overview of the entire system in Section IV. The design of ABAC using homomorphic encryption is presented in Section V. In Section VI, we introduce the concepts of kABAC and k-Anonymous Quantization. We also describe a greedy algorithm to realize kAQ and show a secure procedure to perform quantization without revealing private information. To demonstrate the viability of our approach, we have tested our system using a large collection of iris patterns. The details of the experiments and the results are presented in Section VII. We conclude the paper and discuss future work in Section VIII.

## II. RELATED WORK

The main contributions of our paper are the introduction of the ABAC system concept and a practical design of such a system using iris biometrics. There are other work that deal with the privacy and security issues in biometric systems but their focus are different from this paper. A privacy-protecting technology called “Cancelable Biometrics” has been proposed in [4]. To protect the security of the raw biometric signals, a cancelable biometric system distorts a

biometric signal using a specially designed non-invertible transform so that similarity comparison can still be performed after distortion. Biometric Encryption (BE) described in [5] possesses all the functionality of Cancelable Biometrics, and is immune against the substitution attack because it outputs a key which is securely bound to a biometric. The BE templates stored in the gallery have been shown to protect both the biometrics themselves and the keys. The stored BE template is also called “helper data”. “Helper data” is also used in [6] to assist in aligning a probe with the template that is available only in the transformed domain and does not reveal any information about the fingerprint.

All the above technologies focus on the security and privacy of the biometric signals in the gallery: instead of storing the original biometric signal, they keep only the transformed and non-invertible feature or helper data extracted from the original signal that do not compromise the security of the system even if they are stolen. In these systems, the identity of the user is always recognized by the system after the biometric matching is performed. To the best of our knowledge, there are no other biometric access systems that can provide access control and yet keep the user anonymous. Though our focus is on user anonymity, our design is complementary to cancelable biometrics and it is conceivable to combine features from both types of systems to achieve both data security and user anonymity.

Anonymity in biometric features like faces is considered in [7]. Face images are obfuscated by a face de-identification algorithm in such a way that any face recognition softwares will not be able to reliably recognize de-identified faces. The model used in [7] is the celebrated  $k$ -anonymity model which states that any pattern matching algorithm cannot differentiate an entry in a large dataset from at least  $k - 1$  other entries [8], [3]. The  $k$ -anonymity model is designed for data disclosure protocols and cannot be used for biometric matching for a number of reasons. First, despite the goal of keeping the user anonymous, it is very important of an ABAC system to verify that a user is indeed in the system. Face de-identification techniques provide no guarantee that only faces in the original database will match the de-identified ones. As such, an imposter may gain access by sending an image that is close to an de-identified face. Second, de-identification techniques group similar faces together to facilitate the public disclosure of the data. This is detrimental to anonymity as face clusters may reveal important identity traits like skin color, facial structure, etc.

Another key difference between anonymity in data disclosure and biometric matching is the

need for secure collaboration between two parties – the biometric server and the user. The formal study of such a problem is Secure Multiparty Computation (SMC). SMC is one of the most active research areas in cryptography and has wide applications in electronic voting, online bidding, keyword search and anonymous routing. While there are no previous work that use SMC for biometric matching, many of the basic components in a BAC system can be made secure under this paradigm. They include inner product [9], [10], polynomial evaluation [11], [12], [13], thresholding [14], [15], [16], median [17], matrix computation [18], [19], logical manipulation [20], k-means clustering [21], [22], decision tree [23], [24], [25] and other classifiers [26], [27], [12], [28] etc. A recent tutorial in SMC for signal processing community can be found in [29].

The main hurdle in applying computationally-secure SMC protocols to biometric matching is their high computational complexity. For example, the classical solution to the thresholding problem<sup>1</sup>, or comparing two private numbers  $a$  and  $b$ , is to use Oblivious Transfer (OT) [30]. OT is a SMC protocol for joint table lookup. The privacy of the function is guaranteed by having the entire table encrypted by a pre-computed set of public keys and transmitted to the other party. The privacy of the selection of the table entry is protected based on obfuscating the correct public key among the dummy ones. Even with recent advances in reducing the computational and communication complexity [31], [17], [32], [33], [34], [13], the large table size, the intensive encryption and decryption operations render OT difficult for pixel or sample-level signal processing operations.

A faster but less general approach is to use Homomorphic Encryption (HE) which preserves certain operations in the encrypted domain [35]. Recently, the homomorphic encryption scheme proposed by IBM and Stanford researcher C. Gentry has generated a great deal of excitement in using HE for encrypted domain processing [36]. He proposed using Ideal Lattices to develop a homomorphic encryption system that can preserve both addition and multiplication operations. This solves an open problem on whether there exists a semantically-secure homomorphic encryption system that can preserve both addition and multiplication. On the other hand, his construction is based on protecting the simplest boolean circuit and its generalization to realistic application is questionable. In an interview, Gentry estimates that performing a Google search

<sup>1</sup>This problem is commonly referred to as the Secure Millionaire Problem in SMC literature.

with encrypted keywords would increase the amount of computing time by about a trillion[37] and even this claim is already challenged by others to be too conservative [38].

More practical homomorphic encryptions such as Paillier cryptosystem can only support addition between two encrypted numbers, but do so over a much larger additive plaintext group, thus providing a wide dynamic range for computation [39]. Furthermore, as illustrated in Section III, multiplication between encrypted numbers can be accomplished by randomization and interaction between parties. Recently, Paillier encryption is being applied in a number of fundamental signal processing building blocks [40] including basic classifiers [27] and Discrete Cosine Transform [41] in encrypted domain. Nevertheless, the public-key encryption and decryption processes in any homomorphic encryption still pose a formidable complexity hurdle to overcome. For example, the fastest thresholding result takes around 5 seconds to compare two 32-bit numbers using a modified Paillier encryption system with a key size of 1024 bits [14]. One of the goals of this paper to utilize homomorphic encryption to construct a realistic biometric matching system that can tradeoff computation complexity with user anonymity in a provably secure fashion.

### III. BACKGROUND

We model any biometric signal  $\mathbf{x} = (x_1, \dots, x_n)^T$  as a  $n$ -dimensional vector from a feature space  $F^n$  where  $F$  is a finite field. We also assume the existence of a commutative distance function  $d : F^n \times F^n \rightarrow \mathfrak{R}^+ \cup \{0\}$  that measures the dissimilarity between two biometric signals. In order for the distance to be computable using the operators in the field, we assume that  $F$  to be a subfield of  $\mathfrak{R}$  so that the components of the constituent vectors will be treated as real numbers in the distance computation. The most commonly used distance is the Euclidean distance:

$$d(\mathbf{x}, \mathbf{y})^2 := \|\mathbf{x} - \mathbf{y}\|_2^2 = \sum_{i=1}^n (x_i - y_i)^2 \quad (1)$$

For the iris patterns used in our experiments,  $F$  is the binary field  $Z_2 = \{0, 1\}$  and  $d(\cdot, \cdot)$  is a modified hamming distance defined below [42]:

$$d_H(\mathbf{x}, \mathbf{y})^2 := \frac{\|(\mathbf{x} \otimes \mathbf{y}) \cap \text{mask}_x \cap \text{mask}_y\|_2^2}{\|\text{mask}_x \cap \text{mask}_y\|_2^2} \quad (2)$$

where  $\otimes$  denotes the XOR operation and  $\cap$  denote the bitwise AND.  $\text{mask}_x$  and  $\text{mask}_y$  are the corresponding mask binary vectors that mask the unusable portion of the irises due to occlusion

by eyelids and eyelash, specular reflections, boundary artifacts of lenses, or poor signal-to-noise ratio. As the mask has substantial variation even among feature vectors captured from the same eye, we assume that the mask vectors do not disclose any identity information.

The special distance function and the high dimension of many feature spaces make them less amenable to statistical analysis. There exist mapping functions that can project the feature space  $F^n$  into a lower dimensional space  $\mathfrak{R}^m$  such that the original distance can be approximated by the distance, usually Euclidean, in  $\mathfrak{R}^m$ . The most well-known technique is Principal Component Analysis (PCA) which is optimal if the original distance is Euclidean [43]. For general distances, mapping functions can be derived by two different approaches – the first approach is Multi-dimensional Scaling (MDS) in which an optimal mapping is derived based on minimizing the differences between the two distances over a finite dataset [44]. The second approach is based on distance relationship with random sets of points and include techniques such as Fastmap [45], Lipschitz Embedding [46] and Local Sensitivity Hashing [47]. In our system, we use both PCA and Fastmap for their low computational complexity and good performance. Here we provide a brief review of the Fastmap procedure and will discuss its secure implementation in Section VI. Fastmap is an iterative procedure in which each step selects two random pivot objects  $\mathbf{x}_A$  and  $\mathbf{x}_B$  and computes the projection  $x'$  for any data point  $\mathbf{x}$  as follows:

$$x' := \frac{d(\mathbf{x}, \mathbf{x}_A)^2 + d(\mathbf{x}_A, \mathbf{x}_B)^2 - d(\mathbf{x}, \mathbf{x}_B)^2}{2d(\mathbf{x}_A, \mathbf{x}_B)} \quad (3)$$

The projection in (3) requires only distance relationships. A new distance is then computed by taking into account the existing projection:

$$d'(\mathbf{x}, \mathbf{y})^2 := d(\mathbf{x}, \mathbf{y})^2 - (x' - y')^2 \quad (4)$$

where  $x'$  and  $y'$  are the projections of  $\mathbf{x}$  and  $\mathbf{y}$  respectively. The same procedure can now be repeated using the new distance  $d'(\cdot, \cdot)$ . It has been demonstrated in [45] that using pivot objects that are far apart, the Euclidean distance in the projected space produces a reasonable approximation of the original distance of many different feature spaces.

Using a dissimilarity metric, we can now define the function of a biometric access control system. It is a computational process that involves two parties: a biometric server (Bob) and a user (Alice). Bob is assumed to have a database of  $M$  biometric signals  $DB = \{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ , where  $\mathbf{x}_i = (x_1^i, \dots, x_n^i)^T$  is the biometric signal of member  $i$ . Alice provides a probe  $\mathbf{q}$  and requests

access from the server. Armed with these notations, we first provide a functional definition of a Biometric Access Control system.

*DEFINITION 1:* A *Biometric Access Control* (BAC) system is a computational protocol between two parties, Bob with a biometric database  $DB$  and Alice with a probe  $\mathbf{q}$ , such that at the end of the protocol, Alice and Bob can jointly compute the following value:

$$y_{BAC} := \begin{cases} 1 & \text{if } d(\mathbf{q}, \mathbf{x}_i)^2 < \epsilon \text{ for some } \mathbf{x}_i \in DB \\ 0 & \text{otherwise.} \end{cases} \quad (5)$$

Adding user anonymity to a BAC system results in the following definition:

*DEFINITION 2:* An *Anonymous BAC* (ABAC) system is a BAC system on  $DB$  and  $\mathbf{q}$  with the following properties at the end of the protocol:

- 1) Except for the value  $y_{BAC}$ , Bob has negligible knowledge about  $\mathbf{q}$ ,  $d(\mathbf{q}, \mathbf{x})$ , and the comparison results between  $d(\mathbf{q}, \mathbf{x})^2$  and  $\epsilon$  for all  $\mathbf{x} \in DB$ .
- 2) Except for the value  $y_{BAC}$ , Alice has negligible knowledge about  $\epsilon$ ,  $\mathbf{x}$ ,  $d(\mathbf{q}, \mathbf{x})$ , and the comparison results between  $d(\mathbf{q}, \mathbf{x})^2$  and  $\epsilon$  for all  $\mathbf{x} \in DB$ .

Like any other computationally secure protocols, “negligible knowledge” used in the above definition should be interpreted as, given the available information to a party, the distribution of all possible values of the private input from the other party is computationally indistinguishable from the uniformly random distribution [48]. The first property in Definition 2 defines the concept of user anonymity, i.e. Bob knows nothing about Alice except whether her probe matches one or more biometric signals in  $DB$ . As it has been demonstrated that even the distance values  $d(\mathbf{q}, \mathbf{x}_i)$  are sufficient for an attacker to recreate  $DB$  [49], the second property is designed to disclose the least amount of information to Alice.

It is impossible to design a secure system without considering the possible adversarial behaviors from both parties. Adversarial behaviors are broadly classified into two types: semi-honest and malicious. A dishonest party is called semi-honest if he follows the protocol faithfully but attempts to find out about others’ private data through the communication. A malicious party, on the other hand, will change private inputs or even disrupt the protocol by premature termination. Making the proposed system robust against a wide range of malicious behaviors is beyond the scope of this paper. Here, we assume Bob to be semi-honest but allow certain malicious behaviors from Alice – we assume that Alice will engage in malicious behaviors only if those

behaviors can increase her chance of gaining access, that is turning  $y_{BAC}$  into 1, from using a purely random probe. This is a restricted model because, for example, Alice will not prematurely terminate before Bob reaches the final step in computing  $y_{BAC}$ . Also, Alice will not randomly modify any private input unless such modification will increase her chance of success.

In Section V, we shall provide an implementation of an ABAC system on iris biometrics that is robust under the above security model. The procedure is based on repeated use of a homomorphic encryption system. An encryption system  $Enc(x)$  is homomorphic with respect to an operation  $f_1(\cdot, \cdot)$  in the plaintext domain if there exists another operator  $f_2(\cdot, \cdot)$  in the ciphertext domain such that:

$$Enc(f_1(x, y)) = f_2(Enc(x), Enc(y)). \quad (6)$$

In our system, we choose the Paillier encryption system as it is homomorphic over a large additive plaintext group and thus providing a wide dynamic range for computation. Given a plaintext number  $x \in Z_N$ , the Paillier encryption process is given as follows:

$$Enc_{pk}(x) = \left[ (1 + N)^x \cdot r^N \bmod N^2 \right] \quad (7)$$

where  $N$  is a product of two equal-length secret primes and  $r$  is a random number in  $Z_N$  to ensure semantic security. The public key  $pk$  consists of only  $N$ . The decryption function  $Dec_{sk}(c)$  with  $c \in Z_{N^2}$  and the secret key  $sk$  being the Euler-phi function  $\phi(N)$  is defined by the following two steps:

- 1) Compute  $\hat{m} = \frac{[(c^{\phi(N) \bmod N^2}) - 1]}{N}$  over the integer field;
- 2)  $Dec_{sk}(c) = \hat{m} \cdot \phi(N)^{-1} \bmod N$

The Paillier system is secure under the decisional composite residuosity assumption and we refer interested readers to [50, ch.11] for details. Paillier is homomorphic over addition in  $Z_N$  and the corresponding function is multiplication over the ciphertext field  $Z_{N^2}$ . We can also carry out multiplication with a known plaintext in the encrypted domain. These properties are summarized below:

$$Enc_{pk}(x + y) = Enc_{pk}(x) \cdot Enc_{pk}(y) \quad (8)$$

$$Enc_{pk}(xy) = Enc_{pk}(x)^y \quad (9)$$

Multiplication with a number to which only the ciphertext is known can also be accomplished with a simple communication protocol. Assume that Bob wants to compute  $Enc_{pk}(xy)$  based on

the ciphertexts  $Enc_{pk}(x)$  and  $Enc_{pk}(y)$ . Alice has the secret key  $sk$  but Bob wants to keep  $x$ ,  $y$  and  $xy$  hidden from Alice.  $MULT(Enc_{pk}(x), Enc_{pk}(y))$  (Protocol 1) is a secure protocol that can accomplish this task. It is secure because Alice can gain no knowledge about  $x$  and  $y$  from the uniformly random  $x - r$  and  $y - s$  where  $r$  and  $s$  are two random numbers generated by Bob, and Bob is never exposed to any plaintext related to  $x$  and  $y$ . The complexities of  $MULT(Enc_{pk}(x), Enc_{pk}(y))$  are three encryptions and seven encrypted-domain operations (multiplication and exponentiation) on Bob side, as well as two decryptions and one encryption on Alice side. The communication costs are three encrypted numbers. The homomorphic properties and this protocol will be used extensively throughout this manuscript.

---

**Protocol 1** Private Multiplication  $MULT(Enc_{pk}(x), Enc_{pk}(y))$

---

**Require:** Bob:  $Enc_{pk}(x), Enc_{pk}(y)$ ; Alice:  $sk$

**Ensure:** Bob computes  $Enc_{pk}(xy)$

- 1) Bob sends  $Enc_{pk}(x-r) = Enc_{pk}(x) \cdot Enc_{pk}(-r)$  and  $Enc_{pk}(y-s) = Enc_{pk}(y) \cdot Enc_{pk}(-s)$  to Alice where  $r$  and  $s$  are uniformly random numbers generated by Bob.
- 2) Alice decrypts  $Enc_{pk}(x-r)$  and  $Enc_{pk}(y-s)$ , computes  $Enc_{pk}[(x-r)(y-s)]$  and send it to Bob.
- 3) Bob computes  $Enc_{pk}(xy)$  in the encrypted domain as follows:

$$\begin{aligned} Enc_{pk}(xy) &= Enc_{pk}[(x-r)(y-s) + xs + yr - rs] \\ &= Enc_{pk}[(x-r)(y-s)] \cdot Enc_{pk}(x)^s \cdot Enc_{pk}(y)^r \cdot Enc_{pk}(-rs) \end{aligned}$$


---

#### IV. SYSTEM OVERVIEW

In this section, we provide an overview of the entire design of our efficient anonymous biometric access control system. Again, we will use Bob and Alice to denote the biometric system owner and the user respectively. The overall framework of our proposed system is shown in Figure 1. There are two main processing components in our systems: the preprocessing step and the matching step. While the matching step is executed for every probe, the preprocessing step is executed only once by Bob to compute a *publicly-available* quantization table based on a process called  $k$ -Anonymous Quantization. The purpose of the public table is that, based on

a joint secure-index selection of the table entry between Alice and Bob, Bob can significantly reduce the scope of the similarity search from the entire database  $DB$  to approximately  $k$  candidates. The  $k$ -Anonymous Quantization guarantees that (1) if there is an entry in Bob's database that matches Alice's probe, this entry must be among these candidates, (2) all the candidates are maximally dissimilar so as to provide the least amount of information about Alice's probe, and (3) the public table discloses no information about Bob's database. The details of the  $k$ -Anonymous Quantization and the secure-index selection will be discussed in Section VI.

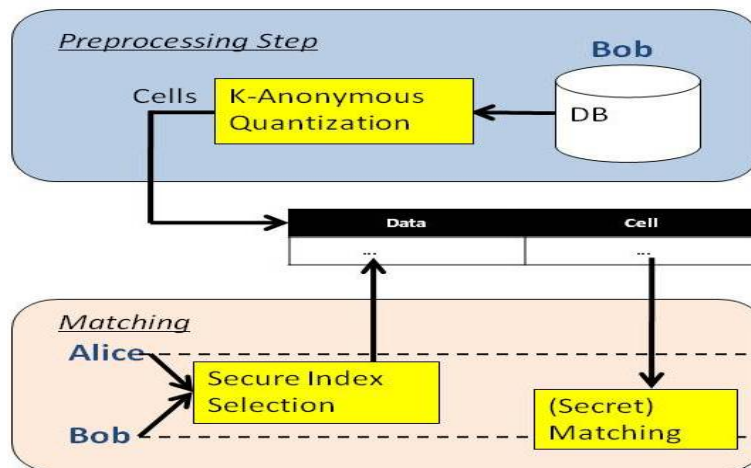


Fig. 1. ABAC System Overview.

After computing the proper quantization cell index from the public table, Bob identifies all the candidates and then engages with Alice in a joint secret matching process to determine if Alice's probe resembles any one of the candidates. This process is conducted in a multi-party computation and communication protocol between Alice and Bob based on Paillier homomorphic encryption. We assume that there is an open network between Bob and Alice that will guarantee message integrity. Since only encrypted content are exchanged, there is no need for any protection against eavesdroppers. For each session, Alice will be responsible for generating the private and public keys for the encryption and sharing the public key with Bob. In other words, a different set of keys will be used for each different user. Furthermore this protocol demands comparable computational capabilities from both parties. Thus it is imperative to use the preprocessing step to reduce the computational complexity of this matching step. As the secret matching utilizes all the fundamental processing blocks for the entire system, we will first explain this component

in the following section.

## V. HOMOMORPHIC ENCRYPTION BASED ABAC

In this section, we describe the implementation of an ABAC system on iris features using homomorphic encryption. The system consists of three main steps: distance computation, bit extraction and secure comparison. Except for the first step of distance computation which is specific towards iris comparison, the remaining two steps and the overall protocol are general enough for other types of biometric features and similarity search. We shall follow a bottom-up approach by first describing individual components and demonstrating their safety before assembling them together as an ABAC system.

### A. Hamming Distance

The modified Hamming distance  $d_H(\mathbf{x}, \mathbf{y})$  described in Equation (2) is used to measure the dissimilarity between iris patterns  $\mathbf{x}$  and  $\mathbf{y}$  which are both 9600 bits long [51]. As the division in Equation (2) may introduce floating point numbers, we focus on the following distance and roll the denominator into the similarity threshold during the later stage of comparison.

$$\widehat{d}_H(\mathbf{x}, \mathbf{y})^2 := \|\ (\mathbf{x} \otimes \mathbf{y}) \cap \text{mask}_{\mathbf{x}} \cap \text{mask}_{\mathbf{y}} \ \|_2^2 \quad (10)$$

DIST (Protocol 2) provides a secure computation of the modified Hamming distances between Alice's probe  $\mathbf{q}$  and Bob's  $DB$ . Alice needs to provide the encryption of individual bits  $\mathbf{q} = (q_1, q_2, \dots, q_n)^T$  and their negation to Bob. Even though Bob can compute the negation in the encryption domain by performing  $Enc_{pk}(\neg q_i) = Enc_{pk}(1 - q_i) = Enc_{pk}(1) \cdot Enc_{pk}(q_i)^{-1}$ , it is computationally more efficient for Alice to compute them in plaintext as demonstrated in Section VII. In step 1a, Bob computes the XOR between each bit of the query and the corresponding bit in each record  $\mathbf{x}_i$ .  $\widehat{d}_H(\mathbf{q}, \mathbf{x}_i)$  can then be computed by summing all the XOR results in the encrypted domain. Bob cannot derive any information about Alice's probe as the operations are all performed in the encrypted domain. Alice does not participate in this protocol at all. The complexity of DIST includes  $O(Mn)$  encrypted-domain operations where  $M$  is the size of  $DB$  and  $n$  is the number of bits for each feature vector.

---

**Protocol 2** Secure computation of distances  $\text{DIST}(DB, \text{Enc}_{pk}(q_j), \text{Enc}_{pk}(\neg q_j))$  for  $j = 1, \dots, n$

---

**Require:** Bob:  $\mathbf{x}_i$  for  $i = 1, \dots, M$ ,  $\text{Enc}_{pk}(q_j)$  and  $\text{Enc}_{pk}(\neg q_j)$  for  $j = 1, \dots, n$

**Ensure:** Bob computes  $\text{Enc}_{pk}[\widehat{d}_H(\mathbf{q}, \mathbf{x}_i)^2]$  for  $i = 1, \dots, M$ .

1) For  $i = 1, \dots, M$ , Bob repeats the following two steps:

a) For  $k = 1, \dots, n$ , compute

$$\text{Enc}_{pk}(q_k \otimes x_k^i) = \begin{cases} \text{Enc}_{pk}(q_k) & \text{if } x_k^i = 0, \\ \text{Enc}_{pk}(\neg q_k) & \text{otherwise} \end{cases}$$

b) Compute

$$\begin{aligned} \text{Enc}_{pk}[\widehat{d}_H(\mathbf{q}, \mathbf{x}_i)^2] &= \text{Enc}_{pk}\left(\sum_{k: [\text{mask}_{\mathbf{q}} \cap \text{mask}_{\mathbf{x}_i}]_i = 1} q_k \otimes x_k^i\right) \\ &= \prod_{k: [\text{mask}_{\mathbf{q}} \cap \text{mask}_{\mathbf{x}_i}]_i = 1} \text{Enc}_{pk}(q_k \otimes x_k^i) \end{aligned}$$


---

## B. Bit Extraction

The next step is to compare the calculated encrypted distance with a plaintext threshold. As comparison cannot be expressed in terms of summation and multiplication of the two numbers, we need to first extract individual bits from the encrypted distance.  $\text{EXTRACT}(\text{Enc}_{pk}(x))$  (Protocol 3) is a secure protocol between Bob and Alice to extract individual encrypted bits  $\text{Enc}_{pk}(x_k)$  for  $k = 1, \dots, l$  from  $\text{Enc}_{pk}(x)$  where  $x$  is a  $l$ -bit number. The idea is for Bob to ask Alice's assistance in decrypting the numbers and extracting the bits. To protect Alice from knowing anything about  $x$ , Bob sends  $\text{Enc}_{pk}(x + r)$  to Alice who then extracts and encrypts individual bits  $\text{Enc}_{pk}[(x + r)_k]$ . Except for the least significant bit (LSB), Bob cannot undo the randomization in  $\text{Enc}_{pk}[(x + r)_k]$  by carrying out an XOR operation with the bits of  $r$  due to the carry bits. To rectify this problem, step 2d in  $\text{EXTRACT}$  zeros out the lower order bits after they have been extracted and stores the intermediate result in  $y$ , thus guaranteeing the absence of any carry bits from the lower order bits during the randomization. Alice cannot learn any information about  $y$  because the bit to be extracted,  $(y + r)_k$ , is uniformly distributed between 0 and 1. Plaintexts obtained by Alice in different iterations are also uncorrelated as a different random number is used by Bob in each iteration. Even though Alice wants to make  $x$  as small as

---

**Protocol 3** Bit Extraction  $\text{EXTRACT}(Enc_{pk}(x))$ 


---

**Require:** Bob:  $Enc_{pk}(x)$  where  $x$  is a  $l$ -bit number; Alice  $sk$ .

**Ensure:** Bob computes  $Enc_{pk}(x_k)$  for  $k = 1, \dots, l$  with  $k = 1$  being the LSB.

- 1) Bob creates a temporary variable  $Enc_{pk}(y) := Enc_{pk}(x)$ .
  - 2) For  $k = 1, \dots, l$ , the following steps are repeated
    - a) Bob generates a random number  $r$  and sends  $Enc_{pk}(y + r)$  to Alice.
    - b) Alice decrypts  $y + r$ , extracts the  $k^{th}$  bit  $(y + r)_k$  and sends  $Enc_{pk}[(y + r)_k]$  back to Bob.
    - c) Bob computes  $Enc_{pk}(x_k) := Enc_{pk}[(y + r)_k \otimes r_k]$ .
    - d) Bob updates  $Enc_{pk}(y) := Enc_{pk}(y - x_k 2^{k-1}) = Enc_{pk}(y) \cdot Enc_{pk}(x_k)^{-2^{k-1}}$
- 

possible to pass the comparison test, there is no advantage of replacing her replies to Bob with any other value. Bob is not able to obtain any information about  $x$  either as all operations are performed in the encrypted domain. Based on the security model introduced in Section III, this protocol is secure. The complexities of EXTRACT are  $l$  encryptions and  $O(l)$  encrypted-domain operation for Bob, as well as  $l$  decryptions and  $l$  encryptions for Alice. The communication costs are  $2l$  encrypted numbers.

### C. Threshold Comparison

Based on the encrypted bit representations of the distances, we can carry out the actual threshold comparison.  $\text{COMPARE}(Enc_{pk}(x_k), y_k$  for  $k = 1, \dots, l)$  (Protocol 4) is based on the secure comparison protocol developed in [52]. Step 2a accumulates the differences between the two numbers starting from the most significant bits. The state variable  $w = 0$  at the  $k^{th}$  step implies that the bits at order  $k$  and higher between  $x$  and  $y$  match perfectly with each other. Step 2b then computes  $Enc_{pk}(c_k)$  where  $c_k = 0$  if and only if  $w = 0$ ,  $x_k = 0$  and  $y_k = 1$ . This implies that  $x < y$ . In other words,  $x < y$  is true if and only if there exists  $c_k = 0$ . In the last step, we invoke the secure multiplication as described in Protocol 1 to combine all  $c_k$  together into  $c$  which is the desired output. Bob gains no knowledge in this protocol as he never handles any plaintext data. The only step that Alice involves in is in the secure multiplication. The adversarial intention of Alice is to make  $c$  zero so as to pass the comparison test. However, the randomization step in

Protocol 1 provides no additional knowledge nor advantage for Alice to change her input. Thus, this protocol is secure. The complexities of COMPARE are  $3l$  encryptions and  $O(l)$  encrypted-domain operations on Bob side, as well as  $2l$  decryptions and  $l$  encryptions on Alice side. The communication costs are  $3l$  encrypted numbers.

---

**Protocol 4** Secure comparison COMPARE( $Enc_{pk}(x_k), y_k$  for  $k = 1, \dots, l$ )

---

**Require** Bob:  $Enc_{pk}(x_k), Enc_{pk}(y_k)$  and  $y_k$  for  $k = 1, \dots, l$ ; Alice:  $sk$

**Ensure** Bob computes  $Enc_{pk}(c)$  such that  $c = 0$  if  $x < y$ .

- 1) Bob sets  $Enc_{pk}(c) := Enc_{pk}(1)$ ,  $Enc_{pk}(w) := Enc_{pk}(0)$ .
  - 2) For  $k = l, \dots, 1$  starting from the MSB, Bob and Alice compute
    - a)  $Enc_{pk}(w) := Enc_{pk}[w + (x_k \otimes y_k)] = Enc_{pk}(w) \cdot Enc_{pk}(x_k \otimes y_k)$
    - b)  $Enc_{pk}(c_k) := Enc_{pk}(x_k - y_k + 1 + w) = Enc_{pk}(x_k) \cdot Enc_{pk}(y_k)^{-1} \cdot Enc_{pk}(1) \cdot Enc_{pk}(w)$
    - c)  $Enc_{pk}(c) := \text{MULT}(Enc_{pk}(c), Enc_{pk}(c_k))$ .
- 

#### D. Overall Algorithm

Protocol 5 defines the overall ABAC system. Steps 1 and 2 show that Alice first sends Bob her public key and the encrypted bits of her probe. Steps 3 and 4 use secure distance computation DIST (Protocol 2) and secure bit extraction EXTRACT (Protocol 3) to compute the encrypted bit representations of all the distances. Steps 4 and 5 then use secure comparison COMPARE (Protocol 4) and accumulate the results into  $Enc_{pk}(u)$  where  $u = 0$  if and only if  $\widehat{d}_H(\mathbf{q}, \mathbf{x}_i)^2 < \epsilon \cdot \|mask_{\mathbf{q}} \cap mask_{\mathbf{x}_i}\|_2^2$  for some  $i$ . To determine if Alice's probe produces a match, Bob cannot simply send Alice  $Enc_{pk}(u)$  for decryption as she will simply return a zero to gain access. Instead, Bob adds a random share  $r$  and sends  $Enc_{pk}(u + r)$  to Alice. The decrypted value  $u + r$  cannot be sent directly to Bob for him to compute  $u$ . Unless  $u = 0$ , the actual value of  $u$  should not be disclosed to Bob in plaintext as it may disclose some information about the distance computations. Instead, we assume the existence of a Collision-Resistant Hash Function HASH to which Bob and Alice share the same key  $pk_H$  [50, ch.4]. Alice and Bob compute  $\text{HASH}_{pk_H}(u+r)$  and  $\text{HASH}_{pk_H}(r)$  respectively. As the hash function is collision resistant, their equality implies that  $u = 0$  and Bob can verify that Alice's probe matches one of the entries in  $DB$  without

---

**Protocol 5** ABAC( $DB, \mathbf{q}$ )

---

**Require:** Bob:  $\mathbf{x}_i, i = 1, \dots, M$  and  $\epsilon$ ; Alice:  $\mathbf{q}$

**Ensure :** Bob computes  $y = 1$  if  $\widehat{d}_H(\mathbf{q}, \mathbf{x}_i)^2 < \epsilon$  for some  $i$  and 0 otherwise

- 1) Alice sends  $pk$  to Bob.
  - 2) Alice computes  $Enc_{pk}(q_j)$  and  $Enc_{pk}(\neg q_j)$  for  $j = 1, \dots, n$  and sends them to Bob.
  - 3) Bob executes  $DIST(DB, Enc_{pk}(q_j), Enc_{pk}(\neg q_j)$  for  $j = 1, \dots, n)$  to obtain  $Enc_{pk}[\widehat{d}_H(\mathbf{q}, \mathbf{x}_i)^2]$  for  $i = 1, \dots, M$ .
  - 4) For  $i = 1, \dots, M$ , Bob and Alice execute  $EXTRACT(Enc_{pk}[\widehat{d}_H(\mathbf{q}, \mathbf{x}_i)^2])$  to obtain the binary representations  $Enc_{pk}[\widehat{d}_H(\mathbf{q}, \mathbf{x}_i)_k^2]$  for  $k = 1, \dots, \lceil \log_2 n \rceil$ .
  - 5) Bob sets  $Enc_{pk}(u) := Enc_{pk}(1)$ .
  - 6) For  $i = 1, \dots, M$ , Bob and Alice computes
    - a)  $Enc_{pk}(c) := COMPARE(Enc_{pk}[\widehat{d}_H(\mathbf{q}, \mathbf{x}_i)_k^2], (\epsilon \| mask_{\mathbf{q}} \cap mask_{\mathbf{x}_i} \|_2^2)_k$  for  $k = 1, \dots, \lceil \log_2 n \rceil$ )
    - b)  $Enc_{pk}(u) := MULT(Enc_{pk}(u), Enc_{pk}(c))$
  - 7) Bob generates a random number  $r$ , computes  $HASH_{pk_H}(r)$  and sends Alice  $Enc_{pk}(u + r)$ .
  - 8) Alice decrypts  $Enc_{pk}(u + r)$ , computes  $HASH_{pk_H}(u + r)$  and sends it back to Bob.
  - 9) Bob sets  $y = 1$  if  $HASH_{pk_H}(r) = HASH_{pk_H}(u + r)$  and 0 otherwise.
- 

knowing the actual value of the probe. Since Alice knows nothing about  $r$ , she cannot cheat by sending a fake hash value. The complexities of Protocol 5 are  $O(M \log_2 n)$  encryptions and  $O(Mn)$  encrypted-domain operations for Bob, as well as  $O(M \log_2 n)$  encryptions and decryptations for Alice. The communication costs are  $O(M \log_2 n)$  encrypted numbers.

## VI. K-ANONYMOUS BAC

In Section V, we show that both the complexities and the communication costs of the ABAC depend linearly on the size of the database, making ABAC difficult to scale to large databases. Inspired by the  $k$ -anonymity model, a simple approach is to tradeoff complexity with privacy by quickly narrowing Alice's query into a small group of  $k$  candidates and then performing the full cryptographic search only on this small group.  $k$  will serve as a parameter to balance between the complexity and the privacy needed by Alice. This is the idea behind the  $k$ -Anonymous

Biometric Access Control system:

*DEFINITION 3:* An  $k$ -Anonymous BAC ( $k$ -ABAC) system is a BAC system on Bob's database  $DB$  and Alice's probe  $\mathbf{q}$  with the following properties at the end of the protocol:

- 1) There exists a subset  $S \subset DB$  with  $|S| \geq k$  such that for all  $\mathbf{x} \in DB \setminus S$ , Bob knows  $d(\mathbf{q}, \mathbf{x})^2 \geq \epsilon$ .
- 2) Except for the value  $y_{BAC}$  as defined in Definition 1, Bob has negligible knowledge about  $\mathbf{q}$  and  $d(\mathbf{q}, \mathbf{x})$ , for all  $\mathbf{x} \in DB$ , as well as the comparison results between  $d(\mathbf{q}, \mathbf{x})^2$  and  $\epsilon$  for all  $\mathbf{x} \in S$ .
- 3) Except for the value  $y_{BAC}$ , Alice has negligible knowledge about  $\epsilon$ ,  $\mathbf{x}$ ,  $d(\mathbf{q}, \mathbf{x})$ , and the comparison results between  $d(\mathbf{q}, \mathbf{x})^2$  and  $\epsilon$  for all  $\mathbf{x} \in DB$ .

The definition of  $k$ -ABAC system is similar to that of ABAC except that Bob can prematurely exclude  $DB \setminus S$  from the comparison. Even though Alice may be aware of such a narrowing process, the  $k$ -ABAC has the same restriction on Alice's knowledge about  $DB$  as the regular ABAC. There are two challenges in designing a  $k$ -ABAC system:

- 1) How do we find  $S$  so that the process will disclose as little information as possible about  $\mathbf{q}$  to Bob?
- 2) How can Alice choose  $S$  that contains the element that is close to  $\mathbf{q}$  without learning anything about  $DB$ ?

Sections VI-A and VI-B describe our approaches in solving these problems in the context of iris matching.

#### A. $k$ -Anonymous Quantization

A direct consequence of Definition 3 is that if there exists a  $\mathbf{x} \in DB$  such that  $d(\mathbf{q}, \mathbf{x})^2 < \epsilon$ ,  $\mathbf{x}$  must be in  $S$ . In order to achieve the goal of complexity reduction, our approach is to devise a static quantization scheme of the feature space  $F^n$  and publish it in a scrambled form so that Alice can select the right group on her own. To explain this scheme, let us start with the definition of a  $\epsilon$ -ball  $k$ -quantization. Define  $B_\epsilon(\mathbf{x})$  or the  $\epsilon$ -ball of  $\mathbf{x}$  to be the smallest subset of  $F^n$  that contains all  $\mathbf{y} \in F^n$  with  $d(\mathbf{y}, \mathbf{x})^2 < \epsilon$ . An  $\epsilon$ -ball  $k$ -quantization of  $DB$  is defined below:

*DEFINITION 4:* An  $\epsilon$ -ball  $k$ -quantization (eBkQ) of  $DB$  is a partition  $\Gamma = \{P_1, \dots, P_N\}$  of  $F^n$  with the following properties:

- 1)  $\bigcup_{i=1}^N P_i = F^n$  and  $P_i \cap P_j = \phi$  for  $i \neq j$ .
- 2) For all  $\mathbf{x} \in DB$ ,  $B_\epsilon(\mathbf{x}) \cap P_j = B_\epsilon(\mathbf{x})$  or  $\phi$  for  $j = 1, \dots, N$ .
- 3)  $|DB \cap P_j| \geq k$  for  $j = 1, \dots, N$ .

Property 1 of Definition 4 ensures that  $\Gamma$  is a partition while property 2 ensures that no  $\epsilon$ -ball centered at a data point straddles two cells. The last property ensures that each cell must at least contain  $k$  elements from  $DB$ . The importance of using an eBkQ  $\Gamma$  is that if  $\Gamma$  is a shared knowledge between Alice and Bob, Alice can select  $P_j \ni \mathbf{q}$  and communicate the *cell index*  $j$  to Bob. Then Bob can compute  $S := DB \cap P_j$  which must contain, if exists, any  $\mathbf{x}$  where  $d(\mathbf{q}, \mathbf{x})^2 < \epsilon$ .

While a typical vector quantization of  $DB$  will satisfy the  $\epsilon$ -ball preserving criteria, the requirement of preserving the anonymity of  $\mathbf{q}$  imposes a very different constraint. Specifically, we would like all the data points in  $S$  to be *maximally dissimilar* so that no common traits can be learned from  $S$ . This leads to our definition of  $k$ -Anonymous Quantization (kAQ):

*DEFINITION 5:* An *optimal  $k$ -anonymous quantization*  $\Gamma^*$  is an eBkQ of  $DB$  that maximizes the following utility function among all possible eBkQ  $\Gamma$ :

$$\min_{P \in \Gamma} \sum_{\mathbf{x}, \mathbf{y} \in P \cap DB} d(\mathbf{x}, \mathbf{y})^2 \quad (11)$$

The utility function (11) can be interpreted as the total dissimilarity of the most homogeneous cell  $P$  in the partition. The utility function also depends on the number of data points in a cell – adding a new point to an existing cell will always increase its utility. Thus finding the partition that maximizes this utility function not only can ensure the minimal amount of dissimilarity within a cell, it also promotes equal distribution of data points among different cells. Given a fixed number of cells, it is important to minimize the variation in the number of data points among different cells so that the computational complexities of the encrypted-domain matching in different cells would be comparable.

It is challenging to solve for the optimal kAQ for the iris matching problem due to the high dimension, 9600 to be exact, and the uncommon distance used. Our first step is to project this high dimensional space into a lower dimensional Euclidean space  $\mathfrak{R}^m$  by using Fastmap followed by PCA. The Fastmap is used to embed the native geometry of the feature space into an Euclidean space while the PCA optimally minimizes the dimension of the resulting space. Even in this lower dimensional space, the structure of a quantization, namely the boundary of

individual cells, can still be difficult to specify. To approximate the boundary with a compact representation, we first use a simple uniform lattice quantization to partition  $\mathfrak{R}^m$  into a rectilinear grid  $\Omega$  consisting of  $L$  bins  $\{B_1, \dots, B_L\}$ . Then, we maximize the utility function (11) but force the cell boundary to be along those of the bins. This turns an optimal partitioning problem in continuous space into a discrete knapsack problem in assigning bins to cells through a mapping function  $f$  to optimize the utility function. The process is described in Figure 2. We denote the resulting approximated  $k$ -quantization as  $\widehat{\Gamma}^*$ .

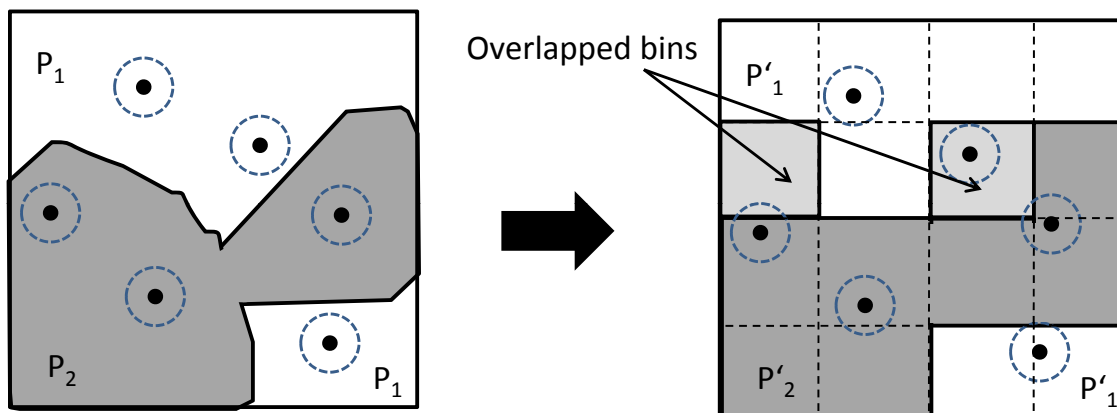


Fig. 2. Approximation of the quantization boundary (left) along the bins (right). The number of bins  $k$  here is 3. There are also two bins that are present in both cells.

As the utility function (11) is based on individual data points, a bin containing multiple  $\epsilon$ -balls may present in multiple cells. As such,  $\widehat{\Gamma}^*$  is no longer a true partition and the mapping function  $f$  is a multi-valued function. A probe falling in these “overlapped” bins will invoke multiple cells, resulting in a larger candidate set  $S$ . Two examples of such overlapped bins are shown in Figure 2. This increases computational complexity and as such, it is important to minimize the amount of overlap. Due to the uneven distribution of data points in the feature space, a global  $\epsilon$  can inflate the size of balls in some area of the feature space resulting in significant overlap problems. In our implementation, we do not use  $\epsilon$  balls but estimate the local similarity structure by using multiple similar feature vectors from each iris, and creating a “bounding box” which is the smallest rectilinear box along the bin boundaries that encloses all the bins containing these similar feature vectors. If any one bin in a bounding box get assigned to cell  $i$ , all the bins in

the bounding box will have an assignment of cell  $i$ .

Protocol 6 (KAQ) describes a greedy algorithm that computes a sub-optimized  $k$ -anonymous quantization mapping function from the data. Step 1 of KAQ sets the number of cells to be the maximum and the protocol will gradually decrease it until each cell has more than  $k$  data points. The initialization steps in 2 and 3 randomly assign a bounding box into each cell. Step 4 identifies the cells that have the minimum utility. Among these cells, steps 5 and 6 identify the cell  $P_{i^*}$  and the bounding box  $BB^*$  which together produce the maximum gain in utility. The bins inside  $BB^*$  are then added to  $P_{i^*}$  and the whole process repeats. This update not only provides a greedy maximization of the overall utility function but also has the tendency to produce an even distribution of data points among different cells. A newly updated cell will have a much lower chance of being updated again as it has a higher utility than others. The final step checks to see if any one cell has less than  $k$  elements and, if yes, restarts the process with fewer target number of cells. For a fixed target number of cells, the complexity of this greedy algorithm is  $O(M^2)$  where  $M$  is the size of  $DB$ . It is important to point out that the output mapping  $f$  only contains entries of bins that belong to at least one bounding box.

### B. Secure Index Selection

Let us first describe how Alice and Bob can jointly compute the projection of Alice's probe  $\mathbf{q}$  into the lower dimensional space formed by Fastmap and PCA. The projection needs to be performed in encrypted domain so that Alice does not reveal anything about her probe and Bob does not reveal any information about his database, the Fastmap pivot points and the PCA basis vectors. Note that the need for encrypted-domain processing does not affect the scalability of our system as the computation complexity depends only on the dimension of the feature space but not on the size of the database.

The Fastmap projection in Equation (3) involves a floating point division. The typical approach of pre-multiplying both sides by the divisor to ensure integer-domain computation does not work. As the Fastmap update Equation (4) needs to square the projection, recursive computation into higher dimensions will lead to a blowup in the dynamic range. To ensure all the computations are performed within in a fixed dynamic range, Alice and Bob need to agree on a pre-defined scaling factor  $\alpha$  and rounding will be performed at each iteration of the Fastmap calculation. Specifically, given the encrypted probe  $Enc_{pk}(\mathbf{q})$ , Bob approximates the first projection  $q'$  in

---

**Protocol 6** Greedy  $k$ -Anonymous Quantization KAQ
 

---

**Require** Bob: Projection of DB into  $\mathfrak{R}^m$  or  $\{P(\mathbf{x}_i)$  for  $i = 1, \dots, M\}$ ; Bin and bounding box structures in  $\Omega$ ;

**Ensure** Bob computes the multi-valued mapping  $f : \Omega \rightarrow \{1, \dots, N\}$  that defines the cell membership of each bin.

- 1) Set the initial number of cells  $N := \lfloor M/k \rfloor$ .
- 2) Let  $L :=$  the list of bounding boxes in  $\Omega$
- 3) Random initialization of cells: for  $i = 1, \dots, N$ ,
  - a) Randomly remove a bounding box  $BB$  from  $L$ .
  - b) Set  $f^{-1}(i) := \{\text{bins in } BB\}$ .
- 4) Identify the collection of cells  $E$  with the lowest utility, i.e.

$$E := \arg \min_{i=1, \dots, N} \sum_{\mathbf{x}, \mathbf{y} \in A_i \cap DB} d(\mathbf{x}, \mathbf{y})^2$$

where  $A_i = \bigcup_{B \in f^{-1}(i)} B$  contains all the bins in cell  $i$ .

- 5) For each cell  $j$  in  $E$ , identify the bounding box  $BB_j^* \in L$  that maximizes the utility of cell  $j$  after adding  $BB_j^*$  to it and denote the resulting utility as  $u_j^*$ , i.e.

$$BB_j^* := \arg \max_{BB \in L} \sum_{\mathbf{x}, \mathbf{y} \in (A_j \cup BB) \cap DB} d(\mathbf{x}, \mathbf{y})^2 \quad (12)$$

$$u_j^* := \sum_{\mathbf{x}, \mathbf{y} \in (A_j \cup BB_j^*) \cap DB} d(\mathbf{x}, \mathbf{y})^2 \quad (13)$$

- 6) Given  $j^* = \arg \max_{j \in E} u_j^*$ , identify the bounding box  $BB^* := BB_{j^*}^*$  and cell  $P_{j^*}$  that give rise to the maximum gain of utility from step 5.
  - 7) Set  $f^{-1}(j^*) := f^{-1}(j^*) \cup \{\text{bins in } BB^*\}$  and remove  $BB^*$  from  $L$ .
  - 8) Go back to Step 4 until  $L$  is empty.
  - 9) For  $i = 1, \dots, N$ , ensure that  $|\bigcup_{B \in f^{-1}(i)} B \cap DB| \geq k$ . If not, set  $N := N - 1$  and go back to step 2.
-

encrypted domain based on the following formula derived from Equation (3):

$$\alpha\tilde{q}' := \text{round}\left(\frac{\alpha}{2ad}\right)\widehat{d}_H(\mathbf{q}, \mathbf{x}_A)^2 + \text{round}\left(\frac{\alpha}{2cd}\right)\widehat{d}_H(\mathbf{x}_A, \mathbf{x}_B)^2 - \text{round}\left(\frac{\alpha}{2bd}\right)\widehat{d}_H(\mathbf{q}, \mathbf{x}_B)^2 \quad (14)$$

where  $a = \|\text{mask}_{\mathbf{q}} \cap \text{mask}_{\mathbf{x}_A}\|_2^2$ ,  $b = \|\text{mask}_{\mathbf{q}} \cap \text{mask}_{\mathbf{x}_B}\|_2^2$ ,  $c = \|\text{mask}_{\mathbf{x}_A} \cap \text{mask}_{\mathbf{x}_B}\|_2^2$  and  $d = d_H(\mathbf{x}_A, \mathbf{x}_B)$ . All the multipliers on the right hand side of (14) are known to Bob in plaintext and the distances can be computed in the encrypted domain using Procedure 2. Since rounding is involved,  $\tilde{q}'$  is just an approximation of  $q'$  as computed with in the original Fastmap formula (3). Based on the computed encrypted values of  $\alpha q'$  from the probe and  $\alpha x'$  from a data point, the update equation (4) is executed as follows:

$$\alpha^2\tilde{d}'_H(\mathbf{x}, \mathbf{q})^2 := \text{round}\left(\frac{\alpha^2}{\|\text{mask}_{\mathbf{x}} \cap \text{mask}_{\mathbf{q}}\|_2^2}\right)\widehat{d}_H(\mathbf{x}, \mathbf{q})^2 - (\alpha\tilde{x}' - \alpha\tilde{q}')^2 \quad (15)$$

Bob again can compute the right hand side of (15) entirely in encryption domain, with the square in the second term computed using Procedure 1. The value  $\widehat{d}'_H(\mathbf{x}, \mathbf{q})^2$  is again approximated due to the rounding of the coefficient. Note that the left hand side has an extra factor of  $\alpha$  which needs to be removed so as to prevent a blowup in the dynamic range. To accomplish that, Bob computes  $\text{Enc}_{pk}(\alpha^2\tilde{d}'_H(\mathbf{x}, \mathbf{q})^2 + r\alpha)$  where  $r$  is a random number, and sends the result to Alice. Alice decrypts it, divides it by  $\alpha$  and round it to obtain  $\text{round}(\alpha\tilde{d}'_H(\mathbf{x}, \mathbf{q})^2) + r$ . Alice encrypts the result and sends it back to Bob who will then removes the random number  $r$ .

Bob can now use the new distances to project the probe along the second pair of pivot objects  $\mathbf{x}_{A'}$  and  $\mathbf{y}_{A'}$  as follows:

$$\alpha^2\tilde{q}'' := \text{round}\left(\frac{\alpha}{2d'}\right)\alpha\tilde{d}'_H(\mathbf{q}, \mathbf{x}_{A'})^2 + \text{round}\left(\frac{\alpha^2}{2}\right) - \text{round}\left(\frac{\alpha}{2d'}\right)\alpha\tilde{d}'_H(\mathbf{q}, \mathbf{x}_{B'})^2 \quad (16)$$

where  $d' = \widehat{d}'_H(\mathbf{x}_{A'}, \mathbf{x}_{B'})^2$  can be computed by Bob in plaintext. The extra factor of  $\alpha$  on the left hand side of (16) can be removed with the help of Alice using a similar approach as previously discussed. As the iteration continues, the deviation of the rounded projection and the original projection will grow as the rounding error accumulates. However, the new distance computed at each iteration absorbs the rounding error from the previous projection. As a result, the distance in the projected space will approach the underlying distance in a similar manner as the original projection.

In the computation of PCA projection, we scale each basis vector with a large enough multiplier not only to absorb the fractional parts of the basis vector but also the scalar  $\alpha$  used

in Fastmap. Let the  $i^{th}$  basis vector of PCA be  $\mathbf{p}_i = \eta(p_1^i, p_2^i, \dots, p_{m_1}^i)^T$  where  $i = 1, \dots, m_2$  with  $m_2$  being the target PCA dimension. The encrypted-domain PCA projection of the Fastmap projection of  $\mathbf{q}$  can be computed as follows:

$$\begin{aligned} Enc_{pk} [P_{pca}(P_{fm}(\mathbf{q}))_i] &:= Enc_{pk} [P_{fm}(\mathbf{q})^T \mathbf{p}_i] = Enc_{pk} \left[ \sum_{j=1}^{m_1} \alpha P_{fm}(\mathbf{q})_j \frac{\eta p_j^i}{\alpha} \right] \\ &= \prod_{j=1}^{m_1} Enc_{pk} [\alpha P_{fm}(\mathbf{q})_j]^{\frac{\eta p_j^i}{\alpha}} \end{aligned} \quad (17)$$

$$\approx \prod_{j=1}^{m_1} Enc_{pk} [\alpha P_{fm}(\mathbf{q})_j]^{\text{round}\left(\frac{\eta p_j^i}{\alpha}\right)} \quad (18)$$

$$(19)$$

The scalar  $\eta$  is selected so that the loss of precision due to rounding is sufficiently small.

The last step of the process is to quantize the projection  $P_{pca}(P_{fm}(\mathbf{q}))$ . We only consider the quantization step size in powers of two so that the quantization process can be performed in the encrypted domain: first, we use the secure bit extraction routine `EXTRACT` to compute the binary representation of  $Enc_{pk} [P_{pca}(P_{fm}(\mathbf{q}))]$ . Then, we drop the lower order bits based on the chosen step-size. The resulting bits are recombined to form the binary representation to the encrypted bin index  $Enc_{pk}(B)$ .

In order to obtain the cell index  $f(B)$ , we need an additional cryptographic tool: a homomorphic collision-resistant hash function  $h_{pk_h}(\cdot)$  with the following homomorphic property [53], [54]:

$$h_{pk_h}(x + y) = h_{pk_h}(x) \cdot h_{pk_h}(y) \quad (20)$$

Our implementation is based on [53]. Bob generates both the public key  $pk_h$  and the secret key for this hash function, and shares the public key with Alice. Instead of directly publishing the mapping  $f(\cdot)$  between the bin index and the corresponding cell indices, Bob publishes an obfuscated mapping  $f'(\cdot)$  such that  $f(B) = f'(h_{pk_h}(B))$ . The hash function sufficiently scrambles all the bin indices so that the distribution of Bob's data among all the bins classified in the KAQ algorithm is disguised as random sampling in the range of the hash function. To prevent Alice from launching a dictionary attack on the table, the length of the bin index must be large enough. This can be accomplished, for example, by padding random projections of the query to make the bin index longer. The cell indices will be published without any obfuscation – little information

is leaked through them as it is shared knowledge between Alice and Bob that there are roughly  $N/k$  distinct cell indices, each of them occurring around  $k$  times.

The reason we need the homomorphic property (20) is to help Alice in computing  $h_{pk_h}(B)$ . After Bob finishes the computation of  $Enc_{pk}(B)$ , he picks a random  $r$ , computes  $h_{pk_h}(r)$  and  $Enc_{pk}(B-r)$  and sends them to Alice. Alice then decrypts  $Enc_{pk}(B-r)$ , computes  $h_{pk_h}(B-r)$  and uses the homomorphic property to compute  $h_{pk_h}(B) = h_{pk_h}(B-r) \cdot h_{pk_h}(r)$ . After that, Alice performs a table lookup to find  $f'(h_{pk_h}(B)) = f(B)$ . If there are multiple cell indices in  $f(B)$ , Alice should not send all of them to Bob because he may use this information to significantly reduce the possible choices of  $B$  as overlapped bins are rare. Instead, Alice should send one cell index first. Then, she re-encrypts her probe and reruns the entire dimension reduction and index selection process as if she was a different user. The same  $f(B)$  will be computed and Alice sends Bob the second index. The whole process is repeated until all the cell indices in  $f(B)$  are exhausted or a match occurs.

SELECT (Protocol 7) summarizes the above process on how Bob can identify the cell to which  $q$  belongs. As for the security of Protocol 7, steps 1 through 4 are processing in encrypted domain and thus reveal no secrets to either parties. Steps 5 and 6 allow Bob to identify the cell indices to which  $q$  belongs. As we assume Bob to be semi-honest, Bob will not deviate from the protocol by adding any identifiable information to the public table  $f'(\cdot)$ . Alice has no incentive to deviate from this protocol as a wrong cell index will erase any chance of success in the subsequent encrypted-domain matching with the elements in the cell. The complexities of Protocol 7 are  $O(m_1 m_2 + m_2 l)$  on Bob side and  $O(m_2 l)$  on Alice side, where  $m_1$  is the Fastmap dimension,  $m_2$  is the PCA dimension and  $l$  is the bit length of the scaled PCA coordinates. The communication costs are  $O(m_1 + m_2 l)$  encrypted numbers.

## VII. EXPERIMENTS AND DISCUSSIONS

For our experiments, we use the CASIA Iris database from the Chinese Academy of Sciences Institute of Automation (CASIA) [55], a common benchmark for evaluating the performance of iris recognition systems. For the iris feature extraction, we use the MATLAB code from [51] to generate both the iris feature vectors and the masks. Each iris feature vector is 9600 bit long. The similarity threshold  $\epsilon$  is set to be 0.35. We select 1,948 samples from CASIA based on the following criteria: the distances are smaller than 0.35 between any two samples from

---

**Protocol 7** Secure Cell Index Selection SELECT
 

---

**Require** Alice: Probe  $\mathbf{q}$ ; Bob: Fastmap pivot objects, PCA basis, and quantization step-size in PCA space,  $\{2^{q_i}$  for  $i = 1, \dots, m_2\}$ ; Public: Scrambled Mapping  $\tilde{f}$ , Deterministic homomorphic cipher with unknown secret key  $Enc_{pk^*, r^*}(\cdot)$

**Ensure** Bob gets  $f(B)$  where  $B \in \Omega$  contains  $\mathbf{q}$

- 1) Alice and Bob computes  $Enc_{pk} [P_{pca}(P_{fm}(\mathbf{q}))_i]$  for  $i = 1, \dots, m_2$ .
  - 2) Bob creates an empty list  $G := \phi$ .
  - 3) Quantization of the projection: for  $i = 1, \dots, m_2$ ,
    - a) Bob and Alice execute  $R := \text{EXTRACT}[Enc_{pk}(P_{pca}(P_{fm}(\mathbf{q}))_i)]$  to get the encrypted binary representation of the  $i^{\text{th}}$  dimension of the projection of  $\mathbf{q}$ .
    - b) Bob discards  $q_i$  lower order encrypted bits from  $R$  and add the remaining bits to the set  $G$ .
  - 4) Bob recombines individual encrypted bits in  $G$  to create a single encrypted  $Enc_{pk}(B)$ .
  - 5) Bob generates a random number  $r$ , compute and sends Alice  $Enc_{pk}(B - r)$  and  $h_{pk_h}(r)$ .
  - 6) Alice decrypts  $Enc_{pk}(B - r)$ , computes  $h_{pk_h}(B) = h_{pk_h}(B - r) \cdot h_{pk_h}(r)$  and uses it look up the cell indices  $f(B) = f'(h_{pk_h}(B))$ .
  - 7) If  $f(B)$  has multiple cell indices, Alice will send the first one to Bob, wait for a random amount of time, re-execute this entire procedure, and sends the second cell index. The process is repeated until all cell indices in  $f(B)$  are exhausted or a match occurs.
- 

the same eye, and larger than 0.40 between any two samples from different eyes. Furthermore, each eye contains at least six good samples and one sample is set aside for testing. A total of 160 individuals are included in our dataset. Our Paillier implementation is based on the Paillier Library developed by J. Bethencourt [56]. The key length of the Paillier cipher is set to be 1024 bit which results in 2048-bit ciphertexts.

### A. Encrypted Domain Processing

In this subsection, we summarize the complexity and communication costs of various encrypted-domain processes discussed in this manuscript. The communication cost is measured based on total amount of information exchanged between Bob and Alice without any overhead from

the network stack. The computation time excludes networking time and is computed based on averaging 100 trials. All of them are implemented in C language on a Linux machine with a 2.4 GHz AMD Athlon 64 CPU and 2 GB memory. Table I summarizes the results. Encrypted-domain addition and multiplication with plaintext are relatively lightweight, except when the plaintext multiplier is negative (i.e. a large positive number in modular arithmetic). Multiplication between two encrypted numbers (MULT) takes the longest and requires information exchange between Bob and Alice. Hamming distance (DIST) is fast as there are no encryption or decryption. Bit extraction (EXTRACT) takes longer and threshold comparison (COMPARE) takes the longest due to the repeated use of negative numbers, encryption and decryption processes. The long computation time for Query preparation is primarily due the high dimension of the iris feature. The overall computation of an ABAC system consists of a fixed setup time of query preparation followed by the time taken for the remaining steps scaled by the size of the database. For a database of 10,000 iris, our ABAC system is estimated to take 41,490 seconds or 11.5 hours and 120 MBytes of network bandwidth. On the other hand, in a  $k$ -anonymous ABAC system, the fixed setup time are the Query Preparation and the SELECT process. The matching complexity depends only on  $k$  but not on the size of the database, except for the rare cases in which the probe falls into an overlapped bin. We shall study the effect of the quantization on the number of overlapped bins in details in Section VII-B. Apart from these exceptions, for the same database of 10,000 iris patterns using a  $k$ -ABAC system with  $k = 50$ , the time required is only 650 seconds and the bandwidth is 1.3 MBytes.

### *B. $k$ -Anonymous Quantization*

In the  $k$ -ABAC system, we first use Fastmap to reduce the original 9600-bits iris code into 100-dimension Euclidean space. Then we use PCA again to further reduce the dimension. Two PCA dimensions, 10 and 20, are tested in our experiments. These steps were performed on a machine running Windows XP Pro. with 3.4 GHz Intel Pentium 4 CPU and 2 GB of RAM. The run time for Fastmap and PCA are 36.24 and 0.274 seconds. There is a loss in performance in each step of projection as the distances cannot be represented as accurately. The plots of False Accept Rates (FAR) versus False Reject Rate (FRR) for the original space and the two projected cases are shown in Figure 3. The performance clearly declines as the dimension decreases from 20 to 10. The consequence of dimension reduction is that the similarity structure cannot be

TABLE I  
TIME AND COMMUNICATION COMPLEXITIES OF ENCRYPTED-DOMAIN PROCESSING

Process	Bob's Time in sec.	Alice's Time in sec.	Communication (Kbits)
Encryption $Enc_{pk}(x)$	$17.3 \times 10^{-3}$	-	-
Decryption $Dec_{sk}(c)$	$12.8 \times 10^{-3}$	-	-
Addition $Enc_{pk}(x) \cdot Enc_{pk}(y)$	$13 \times 10^{-6}$	-	-
Multiplication $Enc_{pk}(x)^y, y \geq 0$	$0.143 \times 10^{-3}$	-	-
Multiplication $Enc_{pk}(x)^y, y < 0$	$30.1 \times 10^{-3}$	-	-
MULT	$47.9 \times 10^{-3}$	$43.0 \times 10^{-3}$	3
DIST <sup>a</sup>	$98 \times 10^{-3}$	-	-
EXTRACT <sup>b</sup>	0.845	0.421	56
COMPARE <sup>b</sup>	2.06	0.602	42
Query Preparation (Step 2 in ABAC)	-	290	-
Remaining steps in ABAC <sup>a</sup>	3.05	1.07	98
SELECT <sup>c</sup>	2149.842	3.455	5522

<sup>a</sup> Average running time for each entry in *DB* amortized over 100 entries, with the dimension of each entry equal to 9600.

<sup>b</sup> 14 bits operand are used as they are sufficient for the Hamming distance.

<sup>c</sup> Fastmap dimension  $m_1 = 100$ ; PCA dimension  $m_2 = 20$  and  $l = 64$ .

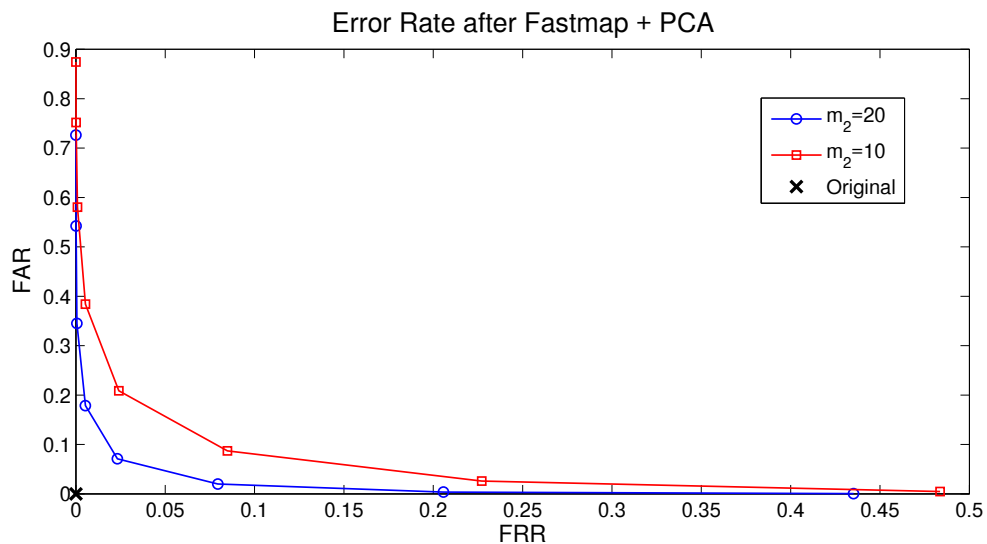


Fig. 3. FRR vs. FAR for using (a) the original feature space, (b) 100d Fastmap and then  $m_2 = 20$  dimensional PCA, and (c) 100d Fastmap and then  $m_2 = 2$  dimensional PCA

well approximated in low dimensions. In defining the  $k$ -Anonymous quantization, we rely on a uniform quantization grid and similarity within a single iris is estimated based on a bounding box of similar features. If the similarity structure is poorly represented, bounding boxes begin to overlap. Probe falling in overlapped areas may need to invoke multiple cells, and thus increase the computational complexities. Figure 4 show the histogram of the fraction of bins that overlap different numbers of bounding boxes. For  $m_2 = 20$ , 88% of the bins are contained in only one bounding box and 96% in at most two bounding boxes. When the dimension is reduced to  $m_2 = 10$ , these numbers reduces to 55% and 76%. Even though overlapped bins are not necessarily classified into different cells by the KAQ algorithm, their total number serve as the upper bound of bins with multiple cell affiliations.

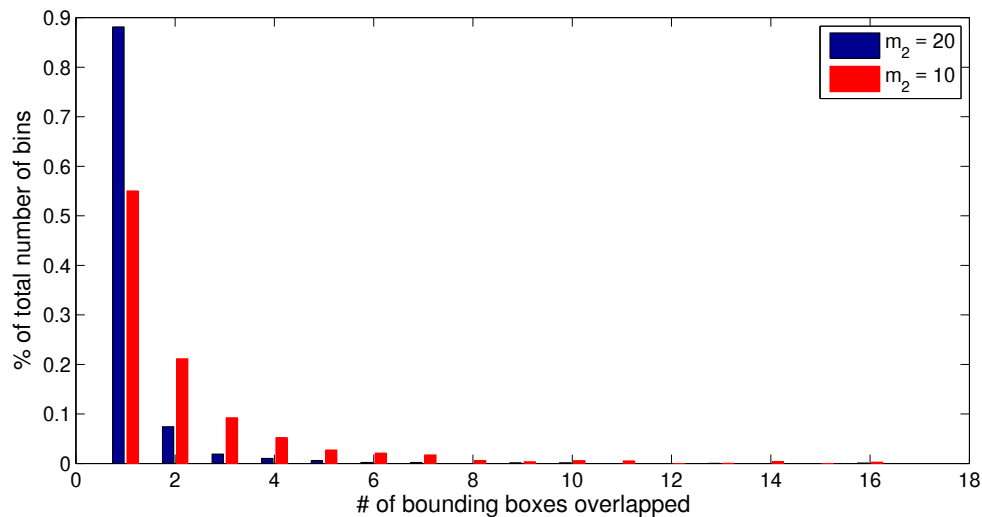


Fig. 4. Histogram of overlapped bins

Next, we consider the performance of KAQ. This algorithm, programmed in C language, was run on a machine running Windows XP Pro. with 2.0 GHz AMD Athlon 64 CPU and 1 GB of RAM. The execution time is a function of the size of the database and takes less than 2 ms to complete regardless of the parameters we used. We have tested the algorithms for various values of  $k$  and for  $m_2 = 10$  and 20 dimensions. Table II summarizes the outputs of the KAQ algorithm at  $m_2 = 20$ . The first column shows the input parameter  $k$ . The second column shows the average and standard deviation of the number of data points in each cell.  $k$  is the lower

TABLE II  
OUTPUT STATISTICS OF THE KAQ ALGORITHM AT  $m_2 = 20$

k	Cell size	Utility	Cell Utility	Complexity
100	106.5±4.0	73856	80262±4885	160.2±146
120	127.2±4.7	106438	115881±5532	189.5±165
150	157.7±5.8	174855	179855±5818	232.1±191
200	207.9±5.1	311756	315252±3016	300.1±226
300	303.5±5.0	673085	679503±8149	423.7±275

bound of the cell size and KAQ manages to produce consistent cell sizes with small variance. The third column shows the utility function as defined in (11) which measures the minimum level of privacy among all the cells. The fourth column considers the average utility function and its standard deviation over all the cells. Again, the standard deviations are generally very small demonstrating the consistency across different cells. The utility increases with  $k$  as the bigger the  $k$  is, the more data points are grouped into the same cell. On the other hand, neither the cell size nor  $k$  are reliable metrics of complexity as they do not take the overlapping among cells into consideration. To provide a more realistic measure, we hold back one data point per individual iris during the quantization construction and use them to test the true complexity. Specifically, we measure complexity based on the actual number of data points in the union of cells that contains the testing probe. The results are tabulated in the last column. The complexity number will be larger than the cell size if the probe falls into a bin that overlaps more than one cell and the number of data points will at least double. The quantized increase in the number of cells accounts for the large standard deviation. In general, the complexity is roughly 1.5 times that of the average cell size.

Table III summarizes the results for KAQ  $m_2 = 10$ . While showing a similar trend as Table II, there are a number of major differences. All the measurements show a much higher level of noise as compared with the previous experiments. This is due to the significant amount of overlapping among bounding boxes. Thus, even when the KAQ algorithm tries to evenly spread the data points, the overlapping forces bounding boxes to be in many cells at the same time. As a consequence, the complexity numbers are much higher than those from KAQ at  $m_2 = 20$ . The utility numbers also decrease from before as the distance measurements are not as well

k	Cell size	Utility	Cell Utility	Complexity
100	153.9±52	44074	87421±67533	567.7±354
120	162.4±49	50965	95583±65760	582.9±355
150	182.5±41	79450	118268±59472	635.3±377
200	224.1±26	145441	176631±42509	724.2±404
300	315.3±7.9	332649	358721±12955	900.8±436

TABLE III

OUTPUT STATISTICS OF THE KAQ ALGORITHM AT  $m_2 = 10$ 

Cell size	Utility	Cell Utility	Complexity
102.5±46	963.0±764	21620.7±18805	183.0±155
121.8±50	2104.8±1694	29927.7±23258	242.9±230
150.8±57	7732.1±3192	45517.4±34276	275.1±237
196.9±65	11747.6±5714	76586.2±48475	327.4±259
285.2±71	50150.8±17737	156620.9±71238	447.3±335

TABLE IV

OUTPUT STATISTICS OF THE RANDOM ALGORITHM AT  $m_2 = 20$ 

preserved.

As there are no comparable quantization schemes in the literature for maximizing privacy, we have chosen, as a reference scheme, random cell assignment for each bounding box at a target number of cells. We call this scheme RANDOM and it is a sensible choice for ensuring individuals with similar iris features to be grouped at a random manner. The testing methodology is that we would first run the KAQ algorithm approach for a specific  $k$ , and then use the same number of cells for RANDOM. Ten random trials of RANDOM are run at each operating point. The results for  $m_2 = 20$  are summarized in Table IV. As expected, RANDOM shows a significant drop in utility as no explicit optimization mechanism is used. The complexity numbers are comparable to those of KAQ as they are mostly a function of the geometry of the data distribution which dictates the overlapping of the bounding boxes.

We finally present the idea of trading off complexities with privacy, as measured by the utility function. We plot the complexity versus utility for all the three schemes in Figure 5. We have

left out the error bars as the standard deviation for the complexity numbers is not meaningful due to the quantized effect of cell increase. This figure demonstrates that the KAQ algorithm provide a good level of privacy protection as the curves for both dimension reside on the high end of utility. While KAQ at  $m_2 = 10$  does not scale well when a high level of privacy is needed, KAQ at  $m_2 = 20$  stays relatively linear. RANDOM is not able to offer much privacy protection.

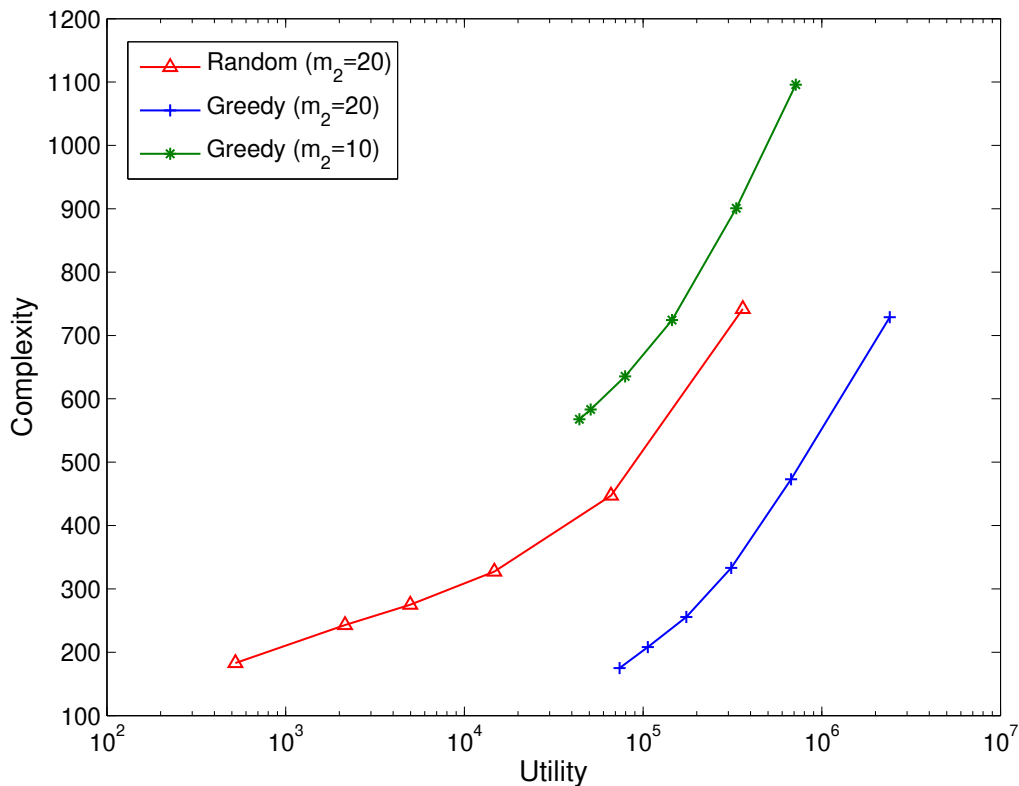


Fig. 5. Tradeoff between complexity and utility (privacy)

## VIII. CONCLUSIONS

In this paper, we have proposed a design for the Anonymous Biometric Control System (ABAC) which allows a biometric server to verify the membership status of a user without knowing his/her identity. The system is composed of various secure multiparty protocols including Hamming distance computation, bit extraction, comparison and result aggregation,

all implemented with a homomorphic cipher. To reduce the computational and communication complexities of such a system, we have proposed a framework called the  $k$ -Anonymous ABAC system that tradeoffs privacy and complexity by quantizing the search space into cells, each of which contains at least  $k$  members. Complexity is reduced by restricting the encrypted domain search process to a small number of cells. Privacy is measured by the dissimilarity of the smallest cell. A greedy quantization scheme on a reduced-dimensional space called  $k$ -Anonymous Quantization has been devised to derive the optimal quantization that maximizes privacy. Secure procedures have been proposed to perform the dimensional reduction and cell lookup. Experimental results on a dataset of iris patterns demonstrate the effectiveness of our techniques in terms of balancing privacy and computational costs. We are currently investigating the extension of the proposed systems to handle a broader class of malicious behaviors. Also, we are interested in improving the efficiency of the homomorphic cipher, particularly in the case when small plaintext numbers are used. Another topic under investigation is the scalability of the  $k$ -Anonymous Quantization to a much larger dataset.

## REFERENCES

- [1] A. Jesdanun, *Youtube, Vacom Agree to Mask Viewer Data*. Associated Press, July 2007.
- [2] W. Hassan and L. Logrippo, "Governance policies for privacy access control and their interactions," in *Feature Interactions in Telecommunication and Software Systems VIII*, D. amyot and L. Logrippo, Eds. IOS Press, 2005, pp. 114–130. [Online]. Available: <http://alloy.mit.edu/community/files/GovernancePoliciesforPrivacy.pdf>
- [3] L. Sweeney, "k-anonymity: a model for protecting privacy," in *International Journal on Uncertainty, Fuzziness and Knowledge-based Systems*, 2002, pp. 557–570.
- [4] N. Ratha, J. Connell, and R. Bolle, "Enhancing security and privacy in biometrics-based authentication systems," *IBM Systems Journal*, vol. 40, no. 3, pp. 614–634, 2001.
- [5] S. Hoque, M. Fairhurst, G. Howells, and F. Deravi, "Feasibility of generating biometric encryption keys," *Electronics Letters*, vol. 41, no. 6, pp. 309–311, 2005.
- [6] U. Uludag and A. Jain, "Securing fingerprint template: Fuzzy vault with helper data," *Proceedings of CVPR Workshop on Privacy Research In Vision*, 2006.
- [7] E. N. Newton, L. Sweeney, and B. Main, "Preserving privacy by de-identifying face images," *IEEE transactions on Knowledge and Data Engineering*, vol. 17, no. 2, pp. 232–243, February 2005.
- [8] V. Ciriani, S. D. C. di Vimercati, S. Foresti, and P. Samarati, "k-anonymity," *Secure Data Management in Decentralized Systems*, 2007.
- [9] O. Goldreich, *Foundations of Cryptography: Volume II Basic Applications*. Cambridge, 2004.
- [10] B. Goethals, S. Laur, H. Lipmaa, and T. Mielikainen, "On private scalar product computation for privacy-preserving data mining," in *The 7th Annual International Conference in Information Security and Cryptology (ICISC2004)*, vol. 3506, 2004, pp. 104–120.

- [11] M. Naor and B. Pinkas, "Oblivious polynomial evaluation," *SIAM J. Comput.*, vol. 35, no. 5, pp. 1254–1281, 2006.
- [12] Y.-C. Chang and C.-J. Lu, "Oblivious polynomial evaluation and oblivious neural learning," *Theoretical Computer Science*, vol. 341, pp. 39–54, 2005.
- [13] M. Naor and B. Pinkas, "Oblivious transfer and polynomial evaluation," in *Proc. 31st Annual ACM Symposium on Theory of Computing*, 1999, pp. 554–567.
- [14] I. Damgard, M. Geisler, and M. Kroigard, "Homomorphic encryption and secure comparison," *International Journal of Applied Cryptography*, vol. 1, no. 1, pp. 22–31, 2008.
- [15] M. Fischlin, "A cost-effective pay-per-multiplication comparison method for millionaires," *Lecture Notes in Computer Science*, vol. 2020, pp. 457–472, 2001.
- [16] A. C. Yao, "Protocols for secure computations," in *Proceedings of the 23rd Annual IEEE Symposium on Foundations of computer science*, 1982.
- [17] G. Aggarwal, N. Mishra, and B. Pinkas, "Secure computation of the kth ranked element," in *Proceedings of Advances in Cryptology - EUROCRYPT 2004: International Conference on the Theory and Applications of Cryptographic Techniques*, 2004, pp. 40–55. [Online]. Available: [citeseer.ist.psu.edu/aggarwal04secure.html](http://citeseer.ist.psu.edu/aggarwal04secure.html)
- [18] E. Kiltz, P. Mohassel, E. Weinreb, and M. Franklin, "Secure linear algebra using linearly recurrent sequences," *LECTURE NOTES IN COMPUTER SCIENCE*, vol. 4392, pp. 291–305, 2007.
- [19] R. Cramer and I. Damgaard, "Secure distributed linear algebra in constant number of rounds," in *Proceedings 21st Annual IACR CRYPTO'01*, ser. LNCS, vol. 2139. Springer-Verlag, 2001, pp. 119–136.
- [20] B. Schoenmakers and P. Tuyls, "Efficient binary conversion for paillier encrypted values," *EUROCRYPT*, vol. 4004, pp. 522–537, 2006.
- [21] G. Jagannathan, K. Pillaipakkamnat, and R. N. Wright, "A new privacy-preserving distributed k-clustering algorithm," *Proceedings of the Sixth SIAM International Conference on Data Mining*, 2006.
- [22] M. C. Doganay, T. B. Pedersen, Y. Saygin, E. Savaş, and A. Levi, "Distributed privacy preserving k-means clustering with additive secret sharing," *Proceedings of the 2008 international workshop on Privacy and anonymity in information society*, pp. 3–11, 2008.
- [23] S. Samet and A. Miri, "Privacy preserving id3 using gini index over horizontally partitioned data," *Computer Systems and Applications, 2008. AICCSA 2008. IEEE/ACS International Conference on*, pp. 645–651, 2008.
- [24] J. Zhan, "Privacy-preserving decision tree classification in horizontal collaboration," *Security of Information and Networks: Proceedings of the First International Conference on Security of Information and Networks (Sin 2007)*, 2008.
- [25] Y. Lindell and B. Pinkas, "Privacy preserving data mining," *Journal of Cryptology*, vol. 15, no. 3, pp. 177–206, 2002.
- [26] J. Vaidya, H. Yu, and X. Jiang, "Privacy-preserving svm classification," *Knowledge and Information Systems*, vol. 14, no. 2, pp. 161–178, 2008.
- [27] C. Orlandi, A. Piva, and M. Barni, "Oblivious neural network computing via homomorphic encryption," *EURASIP Journal on Information Security*, vol. Volume 2007, 2007.
- [28] R. Wright and Z. Yang, "Privacy-preserving bayesian network structure computation on distributed heterogeneous data," *Proceedings of the 2004 ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 713–718, 2004.
- [29] S.-C. Cheung and T. Nguyen, "Secure signal processing between distrusted network terminals," *EURASIP Journal on Information Security*, 2007, <http://www.hindawi.com/GetArticle.aspx?doi=10.1155/2007/51368>.

- [30] M. O. Rabin, "How to exchange secrets by oblivious transfer," Harvar Aiken Computation Laboratory, Tech. Rep. TR-81, 1981.
- [31] D. Boneh, E.-J. Goh, and K. Nissim, "Evaluating 2-dnf formulas on ciphertexts," in *Proceedings of Theory of Cryptography Conference 2005*, ser. LNCS, J. Killian, Ed., vol. 3378. Springer-Verlag, 2005, pp. 325–342.
- [32] M. Naor and B. Pinkas, "Efficient oblivious transfer protocols," in *Proceedings of SODA 2001 (SIAM Symposium on Discrete Algorithms)*, Washington D.C., Jan 2001, pp. 448–457.
- [33] M. Naor and K. Nissim, "Communication complexity and secure function evaluation," *Electronic Colloquium on Computational Complexity (ECCC)*, vol. 8, no. 062, 2001. [Online]. Available: [citeseer.ist.psu.edu/article/naor01communication.html](http://citeseer.ist.psu.edu/article/naor01communication.html)
- [34] C. Cachin, J. Camenisch, J. Kilian, and J. Muller, "One-round secure computation and secure autonomous mobile agents," in *Automata, Languages and Programming*, 2000, pp. 512–523. [Online]. Available: [citeseer.ist.psu.edu/cachin00oneround.html](http://citeseer.ist.psu.edu/cachin00oneround.html)
- [35] C. Fontaine and F. Galand, "A survey of homomorphic encryption for nonspecialists," *EURASIP Journal on Information Security*, vol. 2007, 2007.
- [36] C. Gentry, "Fully homomorphic encryption using ideal lattices," in *Proceedings of the 41st annual ACM symposium on Theory of Computing*, 2009, pp. 169–179.
- [37] M. Cooney, "Ibm touts encryption innovation," *Computer World*, June 25 2009.
- [38] B. Schneier, "Homomprhic encryption breakthrough," in *Schneier on Security*. [http://www.schneier.com/blog/archives/2009/07/homomorphic\\_enc.html](http://www.schneier.com/blog/archives/2009/07/homomorphic_enc.html), 2009.
- [39] P. Pailler, "Public-key cryptosystems based on composite degree residuosity classes," *Proceedings of International Conference on the Theory and Application of Cryptographic Techniques (EUROCRYPT 99)*, vol. vol. 1592, pp. 223–238, May 1999.
- [40] Z. Erkin, A. Piva, S. Katzenbeisser, R. Lagendijk, J. Shokrollahi, G. Neven, and M. Barni, "Protection and retrieval of encrypted multimedia content: When cryptography meets signal processing," *EURASIP Journal on Information Security*, vol. 2007, 2007.
- [41] T. Bianchi, A. Piva, and M. Barni, "Discrete cosine transform of encrypted images," in *Proceedings of IEEE International Conference on Image Processing*, 2008.
- [42] J. Daugman, "How iris recognition works," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 4, pp. 21–30, Jan. 2004.
- [43] J. Hotelling, "Analysis of a complex of statistical variables into principal components," *J. of Educational Psychology*, vol. 24, pp. 417–441, 1933.
- [44] T. F. Cox and M. A. Cox, *Multidimensional scaling*, 2nd ed. Boca Raton : Chapman & Hall, 2001.
- [45] C. Faloutsos and K.-I. Lin, "Fastmap: a fast algorithm for indexing, data-mining and visualization of traditional and multimedia datasets," in *Proceedings of ACM-SIGMOD*, May 1995, pp. 163–174.
- [46] J. Bourgain, "On lipschitz embedding of finite metric spaces in hilbert space," *Israel Journal of Mathematics*, vol. 52, pp. 46–52, 1985.
- [47] A. Gionis, P. Indyk, and R. Motwani, "Similarity search in high dimneions via hashing," in *Proceedings of the 25th International Conference on Very Large Data Bases (VLDB)*, 1999.
- [48] O. Goldreich, *Foundations of Cryptography: Volume 1, Basic Tools*. Cambridge University Press, 2007.

- [49] P. Mohanty, S. Sarkar, and R. Kasturi, "Privacy and security issues related to match scores," in *Proceedings of Computer Vision and Pattern Recognition Workshop*, June 2006, pp. 162–165.
- [50] J. Katz and Y. Lindell, *Introduction To Modern Cryptography*. Chapman and Hall, 2008.
- [51] L. Masek and P. Kovesi, "Matlab source code for a biometric identification system based on iris patterns," The School of Computer Science and Software Engineering, The University of Western Australia, Tech. Rep., 2003.
- [52] I. Damgard, M. Geisler, and M. Kroigard, "Homomorphic encryption and secure comparison," *International Journal of Applied Cryptography*, vol. 1, no. 1, pp. 22–31, 2008.
- [53] D. Filho and P. Barreto, "Demonstrating data possession and uncheatable data transfer," in *Cryptology ePrint Archive*. Report 2206/150, 2006.
- [54] M. Krohn, M. Freedman, and D. Mazieres, "On-the-fly verification of rateless erasure codes for efficient content distribution," in *Proc. IEEE Symposium on Security and Privacy*, 2004, pp. 226–240.
- [55] T. Tan and Z. Sun, "Casia-irisv3," Chinese Academy of Sciences Institute of Automation, <http://www.cbsr.ia.ac.cn/IrisDatabase.htm>, Tech. Rep., 2005.
- [56] J. Bethencourt, *Paillier Library*, UC Berkeley, <http://acsc.csl.sri.com/libpaillier/>.