

# A Formal Study of the Privacy Concerns in Biometric-based Remote Authentication Schemes

Qiang Tang, Julien Bringer, Hervé Chabanne, David Pointcheval

► To cite this version:

Qiang Tang, Julien Bringer, Hervé Chabanne, David Pointcheval. A Formal Study of the Privacy Concerns in Biometric-based Remote Authentication Schemes. L. Chen and Y. Mu. The 4th Information Security Practice and Experience Conference (ISPEC '08), 2008, Sydney, Australie, Australia. Springer-Verlag, Berlin, 4991, pp.56–70, 2008, Lecture notes in computer science. <inria-00419156>

HAL Id: inria-00419156

<https://hal.inria.fr/inria-00419156>

Submitted on 22 Sep 2009

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A Formal Study of the Privacy Concerns in Biometric-based Remote Authentication Schemes<sup>\*</sup>

Qiang Tang<sup>1\*\*</sup>, Julien Bringer<sup>2</sup>, Hervé Chabanne<sup>2</sup>, and David Pointcheval<sup>3</sup>

<sup>1</sup> DIES, EWI, University of Twente, the Netherlands

<sup>2</sup> Sagem Sécurité

<sup>3</sup> Département d'Informatique, École Normale Supérieure  
45 Rue d'Ulm, 75230 Paris Cedex 05, France

**Abstract** With their increasing popularity in cryptosystems, biometrics have attracted more and more attention from the information security community. However, how to handle the relevant privacy concerns remains to be troublesome. In this paper, we propose a novel security model to formalize the privacy concerns in biometric-based remote authentication schemes. Our security model covers a number of practical privacy concerns such as identity privacy and transaction anonymity, which have not been formally considered in the literature. In addition, we propose a general biometric-based remote authentication scheme and prove its security in our security model.

## 1 Introduction

Privacy has become an important issue in many aspects of our daily life, especially in an era of networking where information access may go far beyond our control. When sensitive information such as biometrics is used, the privacy issues become even more important because corruption of such information may be catastrophic for the relevant applications. In this paper we focus on the issue of handling the privacy concerns in remote biometric-based authentication schemes.

### 1.1 Related work

Biometrics, such as fingerprint and iris, have been used to a higher level of security in order to cope with the increasing demand for reliable and highly-usable information security systems, because they have many advantages over typical cryptographic credentials. For example, biometrics are believed to be unique, unforgettable, non-transferable, and they do not need to be stored. One of the most important application areas is biometric-based authentication schemes, where an authentication is simply a comparison between a reference biometric template and a new template extracted during the authentication process. Note that, depending on the type of biometrics, comparison may mean image matching, binary string matching, etc.

Despite of its advantages, in practice, there are some obstacles in a wide adoption of biometrics.

First, biometrics are only approximately stable over the time, therefore, they cannot be directly integrated into most of the existing systems. To address this issue, error-correction concept is widely used in the literature (e.g. [3,4,8,10,11,18,19,25,29]). Employing this concept, some intermediate information (referred to as helper data in some work) is firstly generated based on a reference biometric template, and later, a newly-extracted template could help to recover the reference template or some relevant information if the distance between the templates is small enough (depending on the type of biometrics). Instead of employing this concept, a number of authors also suggest to compare biometric templates directly (e.g. [1,12,34]). Atallah *et al.* [1] propose a method, in which biometric templates are treated as bit strings and subsequently masked and permuted during the authentication process. Du and Atallah [12,34] investigate a number of biometric comparison scenarios by employing secure multi-party computation techniques. Schoenmakers and Tuyls [27] propose to use homomorphic encryption schemes for biometric authentication schemes by employing multi-party computation techniques.

<sup>\*</sup> This work is partially supported by French ANR RNRT project BACH.

<sup>\*\*</sup> The work was done when the author worked as a postdoc researcher at École Normale Supérieure.

Second, biometrics are usually regarded to be sensitive because they uniquely identify an individual. The sensitivity of biometrics lies in the fact that disclosure of biometrics in a certain application leads to the disclosure of the true identity of the involved users in this application. In addition, if the same type of biometrics of a user is used in two applications, then there is an undeniable link for the user's activities in both applications. Nonetheless, it is worth stressing that biometrics are normally considered to be public information. In [20,28,29,31,33], the authors attempt to enhance privacy protection in biometric authentication schemes, where the privacy means that the compromise of the database will not enable the adversary to recover the biometric template. Ratha, Connell, and Bolle [2,24] introduce the concept of *cancelable biometrics* in an attempt to solve the revocation and privacy issues related to biometric information. Ratha *et al.* [23] intensively elaborate this concept in the case of fingerprint-based authentication systems. Recently, Bringer *et al.* [5,6] propose a number of biometric-based authentication protocols which protect the sensitive relationship between a biometric feature and relevant pseudorandom username.

Practical concerns, security issues, and challenges about biometrics have been intensively discussed in the literature (e.g. [2,17,21,24,26,32]). Tuyls, Skoric, and Kevenaar [30] present a summary of cryptographic techniques for dealing with biometrics.

## 1.2 Motivation and contributions

The stability problem concerned with biometric measurements has been paid pretty much attention and investigated very well at this moment. However, privacy issues concerned with biometrics have not been understood well. With respect to biometric-based authentication schemes, we do not have a general formalization of privacy concerns based on a clear system structure. In practice, privacy may mean much more than the adversary cannot recover the user's biometric template. For instance, a user may also want the relationship between its biometric template and username to remain secret in a service, where the user uses a personalized (pseudorandom) username instead of his true name. This requirement might become much stronger if the user wants to multiple registrations under different usernames at the service provider.

In the rest of this paper, we consider the following scenario for biometric-based authentication schemes: Suppose a human user registers at a service provider to consume some service and would like to authenticate himself to the service provider using his biometric (say, his iris). Typically, the user will choose a personalized username and register his reference biometric information under this username. In order to authenticate himself to the service provider, the user presents his username and some fresh biometric information, and then the service provider will perform a matching between the reference biometric information and the fresh biometric information. The contributions of this paper can be summarized as follows.

First, we propose a new system structure for biometric-based remote authentication schemes. In the new structure, there are four types of components, including human user, sensor client, service provider, and database. There are two motivations for us to assume sensor client and service provider to be independent, which means the service provider does not control the sensor client.

1. One is to protect human users' privacy against a malicious service provider. If a malicious service provider controls the sensor client, then it can easily obtain human users' biometric information and potentially manipulate the information.
2. The other is based on the fact that human users may wish to access the service provider wherever they are. In this case, it is natural to make the assumption that sensor client could be provided by another party which has business agreement with the service provider.

Different from any previous system, the database is assumed to be independent from the service provider and serve as a secure storage for biometric information. The motivations for the detachment are as follows.

1. The first is that a user may not trust a service provider to store his biometric template regardless of the transformation which might be applied to the template.
2. The second is that the service provider’s access to the biometric information can be minimized, so is the database’s access. This structure makes it possible to protect human users’ privacy against a malicious service provider or a malicious database. Under the traditional structure, where the service provider controls the database, we do not see how to achieve our privacy goal<sup>1</sup>.
3. The third is that, in practice, the service provider has avoided the responsibility for storing biometric templates. As data breaches for service providers are reported more and more frequently nowadays, the need for the separation becomes stronger and stronger.

With respect to the new structure, we formalize the following attributes related to privacy concerns which have not been formally considered in the literature.

- The security for private relationship between personalized username and biometric template is defined to be an attribute *identity privacy*.
- The security for user’s transaction statistics is defined to be an attribute *transaction anonymity*.

Note that, for non biometric-based (authentication) schemes, the requirement of identity privacy might not be as significant as in our case because cryptographic credentials are not bound to an individual permanently.

Second, we propose a general biometric-based remote authentication scheme by employing a Private Information Retrieval (PIR) protocol [7] (described in the Appendix A) and the ElGamal public-key encryption scheme [13] (described in the Appendix B). The security of the scheme is based on the semantic security of ElGamal, namely the DDH assumption. Instead of ElGamal, other homomorphic encryption schemes can also be used for the same purpose but the computational load will stay in a similar level. Our proposal is not focused on a specific biometric, but rather on such type of biometrics that can be represented as binary strings in the Hamming space and authentication can be done through a binary string matching. For example, iris is one type of such biometrics [16]. For other biometrics, how to construct a secure authentication scheme in our security model remains as an open problem.

### 1.3 Organization

The rest of the paper is organized as follows. In Section 2 we provide some preliminary definitions. In Section 3 we provide the security and privacy definitions for biometric-based remote authentication schemes. In Section 4 we present a new biometric-based remote authentication scheme. In Section 5 we provide security analysis for the new scheme in our security model. In Section 6 we conclude the paper.

## 2 Preliminary definitions

### 2.1 The system structure

In the new system structure for biometric-based authentication schemes, we consider four types of components.

- Human user, which uses his biometric to authenticate himself to a service provider.
- Sensor client, which captures the raw biometric data and extracts a biometric template, and communicates with the service provider.
- Service provider, which deals with human user’s authentication request by querying the database.

---

<sup>1</sup> Especially, applying a one-way function to the biometric template will not be enough to achieve our privacy goal.

- Database, which stores biometric information for users, and works as a biometric template matcher by providing the matching service to the service provider.

*Remark 1.* Different from the local authentication environment, sensor client and service provider are assumed to be independent components in our structure. We consider this to be an appropriate assumption in the remote authentication environment, where human users access the service provider through sensor clients, which are not owned by the service provider but have a business agreement with the service provider.

*Remark 2.* In practice, there might be only very few organizations that can be trusted by human users to store their biometric information though they may want to use their biometrics for the authentication purpose at many service providers. Therefore, in practice we suggest an scenario like that of Single Sign-On systems [22], where biometric information for all service providers are centralizedly stored and managed. In addition, in our security model the centralized database won't be a bottleneck in the sense of security.

For the simplicity of description, in the following discussions, we assume  $N$  users  $U_i$  ( $1 \leq i \leq N$ ) register at a service provider  $\mathcal{S}$ , these users authenticate themselves through a sensor client  $\mathcal{C}^2$ , and the database is denoted as  $\mathcal{DB}$ . Moreover, we would expect users to conduct their authentication services at different service providers while registering their biometric templates in the same (trusted) database.

## 2.2 The authentication workflow

Like most existing biometric-based cryptosystems, we also assume that a biometric-based authentication scheme consists of two phases: an enrollment phase and a verification phase.

1. In the enrollment phase, user  $U_i$  registers his reference biometric information, which is computed based on his reference biometric template  $b_i$ , at the database  $\mathcal{DB}$  and his personalized username  $ID_i$  at the service provider  $\mathcal{S}$ . Note that a human user may have multiple registrations at the same service provider.
2. In the verification phase, user  $U_i$  issues an authentication request to the service provider  $\mathcal{S}$  through the sensor client  $\mathcal{C}$ .  $\mathcal{S}$  matches  $U_i$ 's biometric templates with help from the database  $\mathcal{DB}$ .

## 2.3 Assumptions and trust relationships

We make the following assumptions.

1. **Biometric Distribution assumption:** Let  $H$  be the distance function in a metric space (in this paper, we assume it to be Hamming space). Suppose  $b_i$  and  $b_j$  are the reference biometric templates for Alice and Bob, respectively. There is a threshold value  $\lambda$ , the probability that  $H(b_i, b'_j) > \lambda$  is close to 1 and the probability that  $H(b_i, b'_i) \leq \lambda$  is close to 1, where  $b'_i$  and  $b'_j$  are the templates captured for Alice and Bob at any time.
2. **Liveness assumption:** We assume that, with a high probability, the biometric template captured by the sensor is from a live human user. In other words, it is difficult to produce a fake biometric template that can be accepted by the sensor.
3. **Security link assumption:** The communication links between components are protected with confidentiality and integrity. In practice, the security links can be implemented using a standard protocol such as SSL or TLS.

---

<sup>2</sup> In practice, there may be a number of sensor clients for human users to access the service provider, but this simplification will not affect our security result.

The biometric distribution and the liveness assumptions are indispensable for most of biometric-based cryptosystems and they are considered as a prerequisite for the adoption of biometrics. Note that biometrics are public information, additional credentials are always required to establish security links in order to prevent some well-known attacks (e.g. replay attacks). Therefore, the security link assumption is indeed also assumed in most cryptosystems, though it is not as standard as others.

In a biometric-based authentication system, we assume the following trust relationships.

1. Sensor client is always honest and trusted by all other components. By assuming this trust relationship, the liveness assumption is extended from sensor client to service provider in the following sense: when the service provider receives a username and some fresh biometric information, it can confirm with a high probability that the the fresh biometric information is extracted from a human user which has presented the username to the sensor client.
2. With respect to authentication service, service provider is trusted by human users to make the right decision, and database is trusted by human users and the service provider to store and provide the right biometric information. Only an outside adversary may try to impersonate an honest human user.
3. With respect to privacy concerns, both service provider and database are assumed to be malicious which means they may deviate from the protocol specification, but they will not collude. In reality, an outside adversary may also pose threats to the privacy concerns, however, it has no more advantage than a malicious system component.

### 3 Security model for biometric-based authentication

We first describe some conventions for writing probabilistic algorithms and experiments. The notation  $x \stackrel{R}{\leftarrow} S$  means  $x$  is randomly chosen from the set  $S$ . If  $\mathcal{A}$  is a probabilistic algorithm, then  $\mathcal{A}(\text{Alg}; \text{Func})$  is the result of running  $\mathcal{A}$ , which can have any polynomial number of oracle queries to the functionality  $\text{Func}$ , interactively with  $\text{Alg}$  which answers the oracle queries issued by  $\mathcal{A}$ . For the clarity of description, if an algorithm  $\mathcal{A}$  runs in a number of stages then we write  $\mathcal{A} = (\mathcal{A}_1, \mathcal{A}_2, \dots)$ . As a standard practice, the security of a protocol is evaluated by an experiment between an adversary and a challenger, where the challenger simulates the protocol executions and answers the adversary's oracle queries. Without specification, algorithms are always assumed to be polynomial-time and the security parameter is assumed to be  $\ell$ .

Specifically, in our case, there are two functionalities **Enrollment** and **Verification**, where **Enrollment** can be initiated only once to simulate the enrollment phase and **Verification** can be initiated for any user to start an authentication session for any polynomial times. Without loss of generality, if **Verification** is initiated for  $U_i$ , we write  $\text{Verification}(i)$ .

In addition, we have the following definitions for negligible and overwhelming probabilities.

**Definition 3.** The function  $P(\ell) : \mathbb{Z} \rightarrow \mathbb{R}$  is said to be negligible if, for every polynomial  $f(\ell)$ , there exists an integer  $N_f$  such that  $P(\ell) \leq \frac{1}{f(\ell)}$  for all  $\ell \geq N_f$ . If  $P(\ell)$  is negligible, then the probability  $1 - P(\ell)$  is said to be overwhelming.

#### 3.1 Soundness and impersonation resilience

**Definition 4.** A biometric-based authentication scheme is defined to be sound if it satisfies the following two requirements:

1. With an overwhelming probability, the service provider will accept an authentication request in the following case: sensor client sends  $(ID_i, b)$  in an authentication request, where  $H(b, b_i) \leq \lambda$  and  $b_i$  is the reference template registered for  $ID_i$ .
2. With an overwhelming probability, the service provider will reject an authentication request in the following case: sensor client sends  $(ID_i, b)$  in an authentication request, where  $H(b, b_i) > \lambda$  and  $b_i$  is the reference template registered for  $ID_i$ .

If  $b$ , where  $H(b, b_i) \leq \lambda$ , is extracted from a user different from the user registered under  $b_i$ , then we say false accept occurs. Otherwise, if  $b$ , where  $H(b, b_i) > \lambda$ , is extracted from the user registered under  $b_i$ , then we say false reject occurs. From a cryptographic point of view, the false reject rate and the false accept rate may be very high. However, this issue is irrelevant to our privacy concerns, hence, how to handle them is beyond the scope of our paper.

For authentication schemes, impersonation resilience should be the primary goal, nonetheless, under the security link assumption and the liveness assumption, soundness implies impersonation resilience in our case so that we omit the formalization.

### 3.2 Identity privacy

In practice, a malicious service provider or a malicious database may try to probe the relationships between personalized usernames and biometric templates, though they do not need such information in order to make the system work. Informally, the attribute identity privacy means that, for any personalized username, the adversary knows nothing about the corresponding biometric template. It also implies that the adversary cannot find any linkability between registrations in the case that the same human user has multiple registrations at the service provider.

**Definition 5.** A biometric-based authentication scheme achieves identity privacy if  $\mathcal{A} = (\mathcal{A}_1, \mathcal{A}_2)$  has only a negligible advantage in the following game, where the advantage is defined to be  $|\Pr[e' = e] - \frac{1}{2}|$ .

$$\mathbf{Exp}_A^{\text{Identity-Privacy}} \left| \begin{array}{ll} (i, ID_i, b_i^{(0)}, b_i^{(1)}, (ID_j, b_j) (j \neq i)) \leftarrow \mathcal{A}_1(1^\ell) & \\ b_i = b_i^{(e)} \xleftarrow{R} \{b_i^{(0)}, b_i^{(1)}\} & \\ \emptyset \leftarrow \text{Enrollment}(1^\ell) & \\ e' \leftarrow \mathcal{A}_2(\text{Challenger}; \text{Verification}) & \end{array} \right.$$

Note that the symbol  $\emptyset$  means that there is no explicit output (besides the state information) for the adversary. In the experiment, presumably, the adversary  $\mathcal{A}_2$  will obtain the corresponding information<sup>3</sup> from the challenger. The attack game can be informally rephrased as follows:

1. The adversary  $\mathcal{A}_1$  generates  $N$  pairs of username and relevant biometric template, but provides two possible templates  $(b_i^{(0)}, b_i^{(1)})$  for  $ID_i$ .
2. The challenger randomly chooses a template  $b_i^{(e)}$  for the username  $ID_i$ , and simulates the enrollment phase to generate the parameter for the sensor client, the service provider, and the database.
3. The adversary  $\mathcal{A}_2$  can initiate any (polynomial) number of protocol instances for the verification protocol, and terminates by outputting guess  $e'$ .

In this definition (and Definition 6), the adversary can freely choose the username and biometric template pairs for the enrollment phase, therefore, it models the security for any type of biometric regardless of its distribution in practice. It is worth stressing that, if a scheme achieves identity privacy, then neither a malicious service provider or a malicious database (or an outside adversary which has compromised any of them) can recover any registered biometric template.

As to our knowledge, none of the existing biometric-based authentication schemes (including those in Section 1) achieve identity privacy under our definition. Informally, these scheme suffers from the following vulnerability: Suppose that human users use their iris to authenticate themselves to a service provider  $\mathcal{S}$ . If  $\mathcal{S}$  is malicious (or a hacker which has compromised the biometric database of  $\mathcal{S}$ ), then it can easily determine whether a human being, say Alice, has registered.

<sup>3</sup> The information refers to that of the malicious component at the end of the enrollment phase.

### 3.3 Transaction anonymity

Since the database is supposed to store biometric information, therefore, it might obtain some transaction statistics about the service provider and registered human users. Informally, the attribute transaction anonymity means that, for every query issued by the service provider, a malicious database knows nothing about which user is authenticating himself to the service provider.

**Definition 6.** A biometric-based authentication scheme achieves transaction anonymity if an adversary  $\mathcal{A} = (\mathcal{A}_1, \mathcal{A}_2, \mathcal{A}_3)$  has only a negligible advantage in the following game, where the advantage is defined to be  $|\Pr[e' = e] - \frac{1}{2}|$ .

$$\text{Exp}_{\mathcal{A}}^{\text{Transaction-Anonymity}} \left\{ \begin{array}{ll} (ID_j, b_j) (1 \leq j \leq N) & \leftarrow \mathcal{A}_1(1^\ell) \\ \emptyset & \leftarrow \text{Enrollment}(1^\ell) \\ \{i_0, i_1\} & \leftarrow \mathcal{A}_2(\text{Challenger}, \text{Verification}) \\ i_e & \stackrel{R}{\leftarrow} \{i_0, i_1\} \\ \emptyset & \leftarrow \text{Verification}(i_e) \\ e' & \leftarrow \mathcal{A}_3(\text{Challenger}; \text{Verification}) \end{array} \right.$$

As the adversary is a malicious database, presumably the adversary  $\mathcal{A}_2$  will obtain the corresponding information from the challenger. The attack game can be informally rephrased as follows:

1. The adversary  $\mathcal{A}_1$  generates  $N$  pairs of username and relevant biometric template.
2. The challenger simulates the enrollment phase to generate the parameters.
3. The adversary  $\mathcal{A}_2$  can then initiate any (polynomial) number of protocol instances for the verification protocol. At some point,  $\mathcal{A}_2$  chooses two users  $U_{i_0}, U_{i_1}$  and asks the challenger to initiate an instance for the verification protocol.
4. The challenger chooses  $U_{i_e}$  and initiates an instance for the verification protocol.
5. The adversary  $\mathcal{A}_3$  can continue to initiate any number of protocol instances, and terminates by outputting guess  $e'$ .

## 4 A general biometric-based authentication scheme

In this section we describe a general biometric-based authentication scheme, where the biometric template matching can be done through binary string comparison. We first describe the enrollment phase and the verification phase, and then provide some remarks.

### 4.1 The enrollment phase

In the enrollment phase, every component initializes its parameters as follows.

- $\mathcal{C}$  generates a key pair  $(pk_c, sk_c)$  for a signature scheme ( $\text{KeyGen}, \text{Sign}, \text{Verify}$ ) and publishes the public key  $pk_c$ . In addition,  $\mathcal{C}$  implements a  $(\mathcal{M}, m, \tilde{m}, \lambda)$ -secure sketch scheme ( $\text{SS}, \text{Rec}$ ) [11], where  $\mathcal{M}$  is the space of biometric template,  $m$  and  $\tilde{m}$  can be any values, and  $\lambda$  is the threshold value in the biometric distribution assumption described in Section 2.3.
- $\mathcal{DB}$  generates an ElGamal key pair  $(pk_{db}, sk_{db})$ , where  $pk_{db} = (\mathbb{G}_{db}, q_{db}, g_{db}, y_{db})$ ,  $y_{db} = g_{db}^{x_{db}}$ , and  $sk_{db} = x_{db}$ , and publishes  $pk_{db}$ .
- $\mathcal{S}$  generates an ElGamal key pair  $(pk_s, sk_s)$ , where  $pk_s = (\mathbb{G}_s, q_s, g_s, y_s)$ ,  $\mathbb{G}_s = \mathbb{G}_{db}$ ,  $g_s = g_{db}$ ,  $y_s = g_s^{x_s}$ , and  $sk_s = x_s$ , and publishes  $pk_s$ .
- $U_i$  generates his personalized username  $ID_i$  and registers it at the service provider  $\mathcal{S}$ , and registers  $B_i$  at the database  $\mathcal{DB}$ , where  $b_i$  is  $U_i$ 's reference biometric template and

$$\begin{aligned} B_i &= \text{Enc}((g_s)^{ID_s || ID_i || b_i}, pk_s) \\ &= (B_{i1}, B_{i2}) \end{aligned}$$

Note that  $B_i$  has two components since the encryption scheme is ElGamal. In addition,  $U_i$  (publicly) stores a sketch  $sketch_i = \text{SS}(b_i)$ .



## 4.2 The verification phase

If  $U_i$  wants to authenticate himself to the service provider  $\mathcal{S}$  through the sensor client  $\mathcal{C}$ , they perform as follows.

1. The sensor client  $\mathcal{C}$  extracts  $U_i$ 's biometric template  $b_i^*$  and computes the adjusted template  $b'_i = \text{Rec}(b_i^*, \text{sketch}_i)$ . If  $\mathbf{H}(b_i^*, b'_i) \leq \lambda$ ,  $\mathcal{C}$  sends  $(ID_i, M_{i1}, M_{i2}, \sigma_i)$  to the service provider  $\mathcal{S}$ , where

$$\begin{aligned} X_i &= \text{Enc}((g_s)^{ID_s || ID_i || b'_i}, pk_s) \\ &= (X_{i1}, X_{i2}), \end{aligned}$$

$$M_{i1} = \text{Enc}(X_{i1}, pk_{db}), \quad M_{i2} = \text{Enc}(X_{i2}, pk_{db}),$$

$$\sigma_i = \text{Sign}(ID_s || M_{i1} || M_{i2}, sk_c).$$

Otherwise,  $\mathcal{C}$  aborts the operation.

2.  $\mathcal{S}$  first retrieves the index  $i$  for  $ID_i$  and then forwards  $(M_{i1}, M_{i2}, \sigma_i)$  to the database  $\mathcal{DB}$ .
3.  $\mathcal{DB}$  first verifies the signature  $\sigma_i$ . If the verification succeeds,  $\mathcal{DB}$  decrypts  $M_{i1}$  and  $M_{i2}$  to recover  $X_i$ . For every  $1 \leq \ell \leq N$ , the database randomly selects  $s_t \in \mathbb{Z}_{q_s}$  and computes  $R_t = (X_i \odot B_\ell)^{s_t}$ , where, for any integer  $x$  and two ElGamal ciphertexts  $(c_1, c_2)$  and  $(c_3, c_4)$ , the operator  $\odot$  is defined as follows:  $((c_1, c_2) \odot (c_3, c_4))^x = ((\frac{c_1}{c_3})^x, (\frac{c_2}{c_4})^x)$ .
4. The server runs a PIR protocol to retrieve  $R_i$ . If  $\text{Dec}(R_i, sk_s) = 1$ ,  $\mathcal{S}$  accepts the request; otherwise rejects it.

## 4.3 Remarks on the proposed scheme

It is well known that, with ElGamal scheme, we need to encode the plaintext in a certain way in order to obtain semantic security, however, there is no encoding method which will fully preserve the homomorphic property. In our case, we set  $\mathbb{G}_s = \mathbb{G}_{db}$  and  $g_s = g_{db}$ , so that all plaintexts are exponentiations of  $g_s$  and we avoid the encoding problem. The security will not be affected, as we show in next section.

Under the original definition given in [11], a secure sketch scheme is typically used to preserve the entropy of the input and allow the reconstruction of the input in the presence of a certain amount of noise. In our case, we only need the second functionality, namely the secure sketch scheme is used to remove the noise in the fresh biometric template. Therefore, we allow the parameters  $m$  and  $\tilde{m}$  to be any values. The choice of  $\lambda$  depends on both the type of biometric and the underlying application's requirements on false accept and false reject rates.

User  $U_i$  does not need to register any information, either public or private, at the sensor client, though it need to store some public information, namely the secure sketch. The authentication is conducted through an exact equivalence comparison between the reference template and the adjusted fresh template (say, the output from the secure sketch scheme). As a result, we avoid the need to perform approximate biometric matchings on the service provider side and are able to use the underlying cryptographic techniques. This makes the scheme more scalable and flexible than other similar schemes. Compared with the existing remote authentication schemes (e.g. those in [3,4,8]), the proposed scheme demonstrates our concept of detaching biometric information storage from the service provider and shows a way to enhance human users' privacy in practice. In addition, our scheme also demonstrates a method to transform the existing schemes to satisfy our security definition, i.e. using a combination of plaintext equivalence test and PIR.

The computational complexity is dominated by that of the database  $\mathcal{DB}$  which has to perform  $O(N)$  exponentiations, the sensor client needs to perform 6 exponentiations and sign one message for each authentication attempt, while the service provider only needs to decrypt one message (one exponentiation) to make a decision. In addition, there is some computational load in running the PIR protocol. The communication complexity is dominated by the PIR protocol. If it is instantiated to be the single-database PIR protocol of Gentry and Ramzan [14], then the communication complexity between the service provider and the database is  $O(\ell + d)$ , where  $d$  is the bit-length of an ElGamal ciphertext and  $\ell \geq \log N$  is the security parameter.

## 5 Security analysis of the proposed scheme

### 5.1 Soundness and impersonation resilience

From the biometric distribution assumption and the soundness of the secure sketch, it is straightforward to verify that the proposed authentication scheme is sound under Definition 4. In addition,  $U_i$ 's biometric templates  $b_i$  and  $b'_i$  are encoded in the form  $(g_s)^{ID_s || ID_i || b_i}$  and  $(g_s)^{ID_s || ID_i || b'_i}$ . Hence, if the entropy of the adopted biometric is high, then the service provider and the database, even if they collude, cannot recover the biometric templates based on the Discrete Logarithm assumption.

### 5.2 Security proof for identity privacy

In the verification protocol, even if security sketch is adopted, it is not guaranteed that  $b'_i = b_i$ . Therefore, in the security proof, we assume that the difference pattern, i.e. the distribution of  $b'_i - b_i \bmod q$ , is denoted as *pattern<sub>i</sub>*. In fact, the security results are independent from the difference patterns.

**Lemma 7.** *The proposed scheme achieves identity privacy against malicious  $\mathcal{S}$ , based on the semantic security of the ElGamal scheme and the existential unforgeability of the signature scheme.*

*Proof.* If the proposed scheme does not achieve identity privacy against malicious  $\mathcal{S}$ , we construct an algorithm  $\mathcal{A}'$ , which receives a public key  $pk_{challenger}$  from the ElGamal challenger and runs  $\mathcal{A}$  as a subroutine to break the semantic security of the ElGamal scheme. The proof is done through a sequence of games.

**Game<sub>0</sub>:** In this game,  $\mathcal{A}'$  faithfully answers the oracle queries from  $\mathcal{A}$ . Without loss of generality, suppose  $\mathcal{A}$  has advantage  $Adv_0$ .

**Game<sub>1</sub>:** In this game,  $\mathcal{A}'$  faithfully answers the oracle queries by  $\mathcal{A}$ , except that, for any  $1 \leq k \leq N$ , it rejects any message  $(M_{k1}, M_{k2}, \sigma_k)$  sent to  $\mathcal{DB}$  if the message is not generated by itself (in this case  $\mathcal{A}_2$  has forged a signature). Let this event be  $E_1$  and the advantage of  $\mathcal{A}$  be  $Adv_1$ , then we have  $|Adv_0 - Adv_1| \leq \Pr[E_1]$ .

**Game<sub>2</sub>:** In this game,  $\mathcal{A}'$  sets  $pk_{db} = pk_{challenger}$  which is from the challenger and sends  $m_0, m_1$  to the challenger for a challenge, where

$$m_0 = (y_s)^a (g_s)^{ID_s || ID_i || b_i^{(0)}}, \quad m_1 = (y_s)^a (g_s)^{ID_s || ID_i || b_i^{(1)}},$$

and  $1 \leq a \leq q_s$  is randomly generated. Suppose  $\mathcal{A}'$  receives the challenge  $c_b = \text{Enc}(m_b, pk_{db})$ , where  $b$  is the coin toss of the challenger.  $\mathcal{A}'$  answers the Verification queries from  $\mathcal{A}_2$  as follows:

1. For any  $1 \leq k \leq N$ ,  $\mathcal{A}'$  simulates the message sent by  $\mathcal{C}$  as follows:
  - If  $k = i$ ,  $\mathcal{A}$  randomly selects  $r_1 \in q_s$ , samples  $r_2 \in \text{pattern}_i$ , and generates  $(ID_k, M_{k1}, M_{k2}, \sigma_k)$ , where

$$M_{i1} = \text{Enc}((g_s)^{a+r_1}, pk_{db}),$$

$$M_{i2} = c_b \otimes \text{Enc}((y_s)^{r_1} (g_s)^{r_2}, pk_{db}),$$

$$\sigma_i = \text{Sign}(ID_s || M_{i1} || M_{i2}, sk_c),$$

where, for any two ElGamal ciphertexts  $(c_1, c_2)$  and  $(c_3, c_4)$ , the operator  $\otimes$  is defined as follows:  $(c_1, c_2) \otimes (c_3, c_4) = (c_1 c_3, c_2 c_4)$ .

- Otherwise,  $\mathcal{A}'$  generates  $(ID_k, M_{k1}, M_{k2}, \sigma_k)$  by directly following the protocol specification.
2. For any  $1 \leq k \leq N$ , if  $\mathcal{A}_2$  sends  $(M_{k1}, M_{k2}, \sigma_k)$  to  $\mathcal{DB}$ , then  $\mathcal{A}'$  checks whether the message is generated by itself. If not, then it rejects  $\mathcal{A}_2$ 's request; otherwise, it continues by following the protocol specification.

If  $\mathcal{A}_2$  outputs  $e'$ , then  $\mathcal{A}'$  terminates by outputting  $b' = e'$ .

The simulation is faithful with respect to that in  $\text{Game}_1$ , therefore, the advantage  $Adv_2$  of  $\mathcal{A}$  equals to  $Adv_1$ . Since  $\mathcal{A}'$  uses the same coin toss as the ElGamal challenger, i.e.  $e = b$ , then  $\mathcal{A}'$  wins the game against the semantic security of ElGamal scheme with advantage  $Adv_2$ .

To conclude, the advantage relationships, we have  $|Adv_0 - Adv_2| \leq \Pr[E_1]$ , where  $\Pr[E_1]$  is negligible since the signature scheme is existentially unforgeable. As a result, we get a contradiction, and the lemma follows.  $\square$

**Lemma 8.** *The proposed scheme achieves identity privacy against malicious  $\mathcal{DB}$ , based on the semantic security of the ElGamal scheme.*

*Proof.* If the proposed scheme does not achieve identity privacy against malicious  $\mathcal{DB}$ , then we can construct an algorithm  $\mathcal{A}'$ , which receives a public key  $pk_{\text{challenger}}$  from the ElGamal challenger and runs  $\mathcal{A}$  as a subroutine to break the semantic security of the ElGamal scheme.  $\mathcal{A}'$  is defined as follows:

1.  $\mathcal{A}'$  sets  $pk_s = pk_{\text{challenger}}$  and sends  $m_0, m_1$  to the challenger, where

$$m_0 = (g_s)^{ID_s || ID_i || b_i^{(0)}}, m_1 = (g_s)^{ID_s || ID_i || b_i^{(1)}},$$

and obtains a challenge  $c_b = \text{Enc}(m_b, pk_s)$  where  $b$  is the coin toss of the challenger.

2.  $\mathcal{A}'$  sets  $B_i = c_b$  and faithfully answers the oracle queries from  $\mathcal{A}_2$ . Specifically, for any  $1 \leq k \leq N$ , the Verification queries are simulated as follows:

- If  $k = i$ ,  $\mathcal{A}'$  randomly selects  $r_1 \in q_s$ , samples  $r_2 \in \text{pattern}_i$ , and computes  $(M_{k1}, M_{k2}, \sigma_k)$  as follows.

$$M_{i1} = \text{Enc}(B_{i1}(g_s)^{r_1}, pk_{db}),$$

$$M_{i2} = \text{Enc}(B_{i2}(y_s)^{r_1}(g_s)^{r_2}, pk_{db}),$$

$$\sigma_i = \text{Sign}(ID_s || M_{i1} || M_{i2}, sk_c).$$

- Otherwise,  $\mathcal{A}'$  generates  $(M_{k1}, M_{k2}, \sigma_k)$  by directly following the protocol specification.

3. If  $\mathcal{A}_2$  outputs  $e'$ , then  $\mathcal{A}'$  terminates by outputting  $b' = e'$ .

Note that  $\mathcal{A}'$  uses the same coin toss as the ElGamal challenger, i.e.  $e = b$ . The simulation is faithful, and the advantage of  $\mathcal{A}'$  in attacking the ElGamal scheme is equal to the advantage of  $\mathcal{A}$ . As a result, we get a contradiction, and the lemma now follows.  $\square$

### 5.3 Security proof for transaction anonymity

We next show that the proposed scheme achieves transaction anonymity.

**Lemma 9.** *The proposed scheme achieves transaction anonymity against malicious  $\mathcal{DB}$ , based on the semantic security of the ElGamal scheme and the security (user privacy) of the PIR protocol.*

*Proof.* If the proposed scheme does not achieve transaction anonymity against malicious  $\mathcal{DB}$ , then we can construct an algorithm  $\mathcal{A}'$ , which receives a public key  $pk_{\text{challenger}}$  from the ElGamal challenger and runs  $\mathcal{A}$  as a subroutine to break the semantic security of the ElGamal scheme.  $\mathcal{A}'$  is defined as follows:

1. On receiving the output from  $\mathcal{A}_1$ ,  $\mathcal{A}'$  sets  $pk_s = pk_{\text{challenger}}$  and faithfully answers the Verification queries from  $\mathcal{A}_2$ .

2. On receiving  $\{i_0, i_1\}$  from  $\mathcal{A}_2$ ,  $\mathcal{A}'$  sends  $m_0, m_1$  to the ElGamal challenger, where  $r_0$  is sampled from  $pattern_{i_0}$ ,  $r_1$  is sampled from  $pattern_{i_1}$ ,

$$m_0 = (g_s)^{ID_s || ID_i || b_{i_0}} (g_s)^{r_0}, \quad m_1 = (g_s)^{ID_s || ID_i || b_{i_1}} (g_s)^{r_1},$$

and obtains a challenge  $c_b = \text{Enc}(m_b, pk_s)$  where  $b$  is the coin toss of the challenger. Let  $c_b = (\gamma_1, \gamma_2)$ .  $\mathcal{A}'$  sends  $(M_{i_b1}, M_{i_b2}, \sigma_{i_b})$  to  $\mathcal{A}$ , where

$$\begin{aligned} M_{i_b1} &= \text{Enc}(\gamma_1, pk_{db}), \quad M_{i_b2} = \text{Enc}(\gamma_2, pk_{db}), \\ \sigma_i &= \text{Sign}(ID_s || M_{i_b1} || M_{i_b2}, sk_c). \end{aligned}$$

Then  $\mathcal{A}'$  flips a coin  $e$  and runs the PIR protocol to retrieve  $R_{i_e}$ .

3.  $\mathcal{A}'$  faithfully answers the oracle queries by  $\mathcal{A}_3$ . If  $\mathcal{A}_3$  finally outputs  $e'$ , then  $\mathcal{A}'$  terminates by outputting  $b' = e'$ .

Let the event  $e = b$  be  $E_1$ , then  $\Pr[E_1] = \frac{1}{2}$ . If  $E_1$  occurs, then this is a valid attack game for  $\mathcal{A}$  and its advantage is  $Adv = |\Pr[e' = e | E_1] - \frac{1}{2}|$ . It is straightforward to verify that the following equation holds for some negligible  $\epsilon$ ,

$$|\Pr[e' = e | E_1] + \Pr[e' = e | \neg E_1] - 1| = \epsilon \quad (1)$$

otherwise, we can construct an adversary for the PIR protocol. From equation (1), we have the following probability relationships.

$$\begin{aligned} \Pr[b = b'] &= \Pr[E_1] \Pr[e' = e | E_1] + \Pr[\neg E_1] \Pr[e' \neq e | \neg E_1] \\ &= \frac{1}{2} (\Pr[e' = e | E_1] + \Pr[e' \neq e | \neg E_1]) \\ &= \frac{1}{2} + \frac{1}{2} (\Pr[e' = e | E_1] - \Pr[e' = e | \neg E_1]) \\ &\geq \frac{1}{2} + \frac{1}{2} (\Pr[e' = e | E_1] - (1 - \Pr[e' = e | E_1] + \epsilon)) \\ &= \frac{1}{2} + (\Pr[e' = e | E_1] - \frac{1}{2}) - \frac{\epsilon}{2} \end{aligned}$$

$$\begin{aligned} |\Pr[b = b'] - \frac{1}{2}| &= |\frac{1}{2} + (\Pr[e' = e | E_1] - \frac{1}{2}) - \frac{\epsilon}{2} - \frac{1}{2}| \\ &\geq Adv - \frac{\epsilon}{2} \end{aligned}$$

Based on the assumption that ElGamal scheme is semantically secure, then we get a contradiction. The lemma now follows.  $\square$

#### 5.4 Further remarks

In our security analysis, as to an outside adversary, we only considered the case where it has not compromised any system component. If the adversary has compromised the sensor client  $\mathcal{C}$ , then it may impersonate an honest user to the service provider if it obtains this user's biometric template (note that biometrics are public information). This is a common problem for many authentication systems, unless we adopt a tamper-resistant sensor client. If the adversary has compromised the service provider  $\mathcal{S}$  or the database  $\mathcal{DB}$ , then the identity privacy property is still preserved. A possible vulnerability when  $\mathcal{DB}$  is compromised is that it may be able to impersonate any user in the system by impersonating  $\mathcal{DB}$  to the service provider. Again, this is a common problem for most authentication systems, and one possible solution is to adopt a layered security design. For example, tamper-resistant hardware can be used for establishing communication links. Then, even if the adversary has compromised the database, the ciphertexts of biometric templates will not help him to impersonate any honest user.

## 6 Conclusion

In this paper we have proposed a specifically-tailored system structure and security model for biometric-based authentication schemes. In our security model, we describe two privacy properties, namely identity privacy and transaction anonymity, which are believed to be serious concerns because of the uniqueness of biometrics. We have also proposed a general authentication scheme which fulfills the security properties described in our security model. An interesting characteristic of our scheme is that, assuming biometric template and secure sketch to be public, a user does not need to store any private information and register any information at the sensor client. In addition, the security requirements on the secure sketch scheme can be greatly relaxed (entropy preservation is not required). As a further research direction, it is interesting to investigate more efficient solutions in our security model.

## References

1. M. J. Atallah, K. B. Frikken, M. T. Goodrich, and R. Tamassia. Secure biometric authentication for weak computational devices. In *Financial Cryptography*, pages 357–371, 2005.
2. R. M. Bolle, J. H. Connell, and N. K. Ratha. Biometric perils and patches. *Pattern Recognition*, 35(12):2727–2738, 2002.
3. X. Boyen. Reusable cryptographic fuzzy extractors. In V. Atluri, B. Pfitzmann, and P. D. McDaniel, editors, *CCS '04: Proceedings of the 11th ACM conference on Computer and communications security*, pages 82–91. ACM Press, 2004.
4. X. Boyen, Y. Dodis, J. Katz, R. Ostrovsky, and A. Smith. Secure remote authentication using biometric data. In R. Cramer, editor, *Advances in Cryptology - EUROCRYPT 2005*, volume 3494 of *Lecture Notes in Computer Science*, pages 147–163. Springer, 2005.
5. J. Bringer, H. Chabanne, M. Izabachène, D. Pointcheval, Q. Tang, and S. Zimmer. An application of the Goldwasser-Micali cryptosystem to biometric authentication. In J. Pieprzyk, H. Ghodosi, and E. Dawson, editors, *Information Security and Privacy, 12th Australasian Conference, ACISP 2007 Proceedings*, volume 4586 of *Lecture Notes in Computer Science*, pages 96–106. Springer, 2007.
6. J. Bringer, H. Chabanne, D. Pointcheval, and Q. Tang. Extended private information retrieval and its application in biometrics authentications. In *To appear in Proceedings of CANS 2007*, 2007.
7. B. Chor, E. Kushilevitz, O. Goldreich, and M. Sudan. Private information retrieval. *J. ACM*, 45(6):965–981, 1998.
8. G. D. Crescenzo, R. Graveman, R. Ge, and G. Arce. Approximate message authentication and biometric entity authentication. In A. S. Patrick and M. Yung, editors, *Financial Cryptography and Data Security, 9th International Conference*, volume 3570 of *Lecture Notes in Computer Science*, pages 240–254. Springer, 2005.
9. G. D. Crescenzo, T. Malkin, and R. Ostrovsky. Single database private information retrieval implies oblivious transfer. In B. Preneel, editor, *Advances in Cryptology - EUROCRYPT 2000*, volume 1807 of *Lecture Notes in Computer Science*, pages 122–138. Springer, 2000.
10. Y. Dodis, J. Katz, L. Reyzin, and A. Smith. Robust fuzzy extractors and authenticated key agreement from close secrets. In C. Dwork, editor, *Advances in Cryptology - CRYPTO 2006*, volume 4117 of *Lecture Notes in Computer Science*, pages 232–250. Springer, 2006.
11. Y. Dodis, L. Reyzin, and A. Smith. Fuzzy extractors: How to generate strong keys from biometrics and other noisy data. In C. Cachin and J. Camenisch, editors, *Advances in Cryptology - EUROCRYPT 2004*, volume 3027 of *Lecture Notes in Computer Science*, pages 523–540. Springer, 2004.
12. W. Du and M. J. Atallah. Secure multi-party computation problems and their applications: a review and open problems. In *NSPW '01: Proceedings of the 2001 workshop on New security paradigms*, pages 13–22. ACM Press, 2001.
13. T. ElGamal. A public key cryptosystem and a signature scheme based on discrete logarithms. In G. R. Blakley and D. Chaum, editors, *Advances in Cryptology, Proceedings of CRYPTO '84*, volume 196 of *Lecture Notes in Computer Science*, pages 10–18. Springer, 1985.
14. C. Gentry and Z. Ramzan. Single-database private information retrieval with constant communication rate. In L. Caires, G. F. Italiano, L. Monteiro, C. Palamidessi, and M. Yung, editors, *Automata, Languages and Programming, 32nd International Colloquium, ICALP 2005*, volume 3580 of *Lecture Notes in Computer Science*, pages 803–815. Springer, 2005.
15. Y. Gertner, Y. Ishai, E. Kushilevitz, and T. Malkin. Protecting data privacy in private information retrieval schemes. In *Proceedings of the Thirtieth Annual ACM Symposium on the Theory of Computing*, pages 151–160, 1998.
16. F. Hao, R. Anderson, and J. Daugman. Combining crypto with biometrics effectively. *IEEE Transactions on Computers*, 55(9):1081–1088, 2006.
17. J. D. Woodward Jr., N. M. Orlans, and P. T. Higgins. *Biometrics (Paperback)*. McGraw-Hill/OsborneMedia, 2002.
18. A. Juels and M. Sudan. A fuzzy vault scheme. *Des. Codes Cryptography*, 38(2):237–257, 2006.
19. A. Juels and M. Wattenberg. A fuzzy commitment scheme. In *ACM Conference on Computer and Communications Security*, pages 28–36, 1999.

20. J. M. G. Linnartz and P. Tuyls. New shielding functions to enhance privacy and prevent misuse of biometric templates. In J. Kittler and M. S. Nixon, editors, *Audio-and Video-Based Biometric Person Authentication, 4th International Conference*, volume 2688 of *Lecture Notes in Computer Science*, pages 393–402. Springer, 2003.
21. D. Maltoni, D. Maio, A.K. Jain, and S. Prabhakar. *Handbook of Fingerprint Recognition*. Springer, 2003.
22. A. Pashalidis and C. J. Mitchell. A taxonomy of single sign-on systems. In R. Safavi-Naini and J. Seberry, editors, *Information Security and Privacy, 8th Australasian Conference, ACISP 2003*, volume 2727 of *Lecture Notes in Computer Science*, pages 249–264. Springer, 2003.
23. N. Ratha, J. Connell, R. M. Bolle, and S. Chikkerur. Cancelable biometrics: A case study in fingerprints. In *ICPR '06: Proceedings of the 18th International Conference on Pattern Recognition*, pages 370–373. IEEE Computer Society, 2006.
24. N. K. Ratha, J. H. Connell, and R. M. Bolle. Enhancing security and privacy in biometrics-based authentication systems. *IBM Systems Journal*, 40(3):614–634, 2001.
25. R. Safavi-Naini and D. Tonien. Fuzzy universal hashing and approximate authentication. Cryptology ePrint Archive: Report 2005/256, 2005.
26. B. Schneier. Inside risks: the uses and abuses of biometrics. *Commun. ACM*, 42(8):136, 1999.
27. B. Schoenmakers and P. Tuyls. Efficient binary conversion for paillier encrypted values. In *EUROCRYPT*, pages 522–537, 2006.
28. P. Tuyls, A. H. M. Akkermans, T. A. M. Kevenaar, G. Jan Schrijen, A. M. Bazen, and R. N. J. Veldhuis. Practical biometric authentication with template protection. In T. Kanade, A. K. Jain, and N. K. Ratha, editors, *Audio-and Video-Based Biometric Person Authentication, 5th International Conference*, volume 3546 of *Lecture Notes in Computer Science*, pages 436–446. Springer, 2005.
29. P. Tuyls and J. Goseling. Capacity and examples of template-protecting biometric authentication systems. In *ECCV Workshop BioAW*, pages 158–170, 2004.
30. P. Tuyls, B. Skoric, and T. Kevenaar. *Security with Noisy Data*. Springer London, 2008.
31. P. Tuyls, E. Verbitskiy, J. Goseling, and D. Denteneer. Privacy protecting biometric authentication systems: an overview. In *EUSIPCO 2004*, 2004.
32. U. Uludag, S. Pankanti, S. Prabhakar, and A. K. Jain. Biometric cryptosystems: Issues and challenges. *Proceedings of the IEEE*, 92(6):948–960, 2004.
33. E. Verbitskiy, P. Tuyls, D. Denteneer, and J. P. Linnartz. Reliable biometric authentication with privacy protection. In *SPIE Biometric Technology for Human Identification Conf.*, 2004.
34. M. J. Atallah W. Du. Protocols for secure remote database access with approximate matching. Technical report, CERIAS, Purdue University, 2000. CERIAS TR 2000-15.

## Appendix A: An Introduction to PIR

In a PIR protocol, a database contains a bit string  $X = x_1x_2 \cdots x_n$ , where  $x_i \in \{0, 1\}$  for every  $1 \leq i \leq n$ , and a user can run the protocol to retrieve any  $x_i$  from the database (without loss of generality, through a query( $i$ ) query).

In [9,15], a PIR protocol is defined to be secure, if it satisfies the following two requirements.

- soundness: If both the database and the user follow the protocol specification, then, for any query, the user always obtains the bit it wants.
- user privacy: For any constant  $c$ , there exists a security parameter  $k^*$ , for all  $k \geq k^*$ , the following in-equation holds for any  $X \in \{0, 1\}^n$ , any  $1 \leq i, j \leq n$ , and any distinguisher  $D$  implemented by the database:

$$|\Pr[D(X, \text{query}(i)) = 1] - \Pr[D(X, \text{query}(j)) = 1]| \leq \frac{1}{k^c}.$$

Informally the user privacy says that a curious database does not have any information about which bit the user queries. In [9,15], another property, namely data privacy, is also defined for PIR, and we refer the reader to the paper for details.

## Appendix B: An Introduction to ElGamal Scheme

The algorithms (Gen, Enc, Dec) of the ElGamal public key encryption scheme [13] are defined as follows:

1. The key generation algorithm Gen takes a security parameter  $1^\ell$  as input and generates two primes  $p, q$  satisfying  $q|p-1$ . Let  $\mathbb{G}$  be the subgroup of order  $q$  in  $\mathbb{Z}_p^*$ ,  $g$  be a generator of  $\mathbb{G}$ . The private key  $x$  which is randomly chosen from  $\mathbb{Z}_q$ , and the public key is  $y = g^x$ . Let  $\Omega$  be a bijective map from  $\mathbb{Z}_q$  to  $\mathbb{G}$ .

2. The encryption algorithm **Enc** takes a message  $m$  and the public key  $y$  as input, and outputs the ciphertext  $c = (c_1, c_2) = (g^r, y^r \Omega(m))$  where  $r$  is randomly chosen from  $\mathbb{Z}_g^*$ .
3. The decryption algorithm **Dec** takes a ciphertext  $c = (c_1, c_2)$  and the private key  $x$  as input, and outputs the message  $m = \Omega^{-1}(c_1^{-x} c_2)$ .

It is well-known that the ElGamal scheme is semantically secure based on the DDH assumption.