

# Accredited Symmetrically Private Information Retrieval

Mohamed Layouni

School of Computer Science, McGill University,  
3480 University Street, Montreal, H3A 2A7, Quebec, Canada

**Abstract.** With the digitization of society and the continuous migration of services to the electronic world, individuals have lost significant control over their data. In this paper, we consider the problem of protecting personal information subjects. More specifically, we propose a new primitive allowing a data subject to decide when, how, and by whom his data can be accessed, without the database manager learning anything about his identity, at the time the data is retrieved. The proposed solution, which we call *Accredited SPIR*, combines symmetrically private information retrieval and privacy-preserving digital credentials. We present three constructions based on the discrete logarithm and RSA problems. Despite the added privacy safeguards, the extra cost incurred by our constructions is negligible compared to that of the underlying building blocks.

**Keywords:** Symmetrically private information retrieval, anonymous credentials, policy enforcement.

## 1 Introduction

In a transaction-based world, with continuously shrinking resources, access control has always been, and still continue to be a central issue. Oftentimes, to benefit from a service or a resource, one is asked to show his identity, or prove possession of a set of qualifications and privileges. In many cases, this forces individuals into leaving identity trails behind them, which could be used for criminal activities such as unlawful monitoring and identity theft. The data collected from such interactions, although generally rich in personal information, is in most cases stored in databases lying outside the control of the data subject. Various techniques have been proposed in the past to strengthen users' privacy and help protect their personal information. Among these we note privacy preserving digital credentials [Cha85, Bra00, CL02], and symmetrically private information retrieval protocols [GIKM98, CMO00, KO97, AIR01, Lip05].

In a symmetrically private information retrieval (SPIR) system, there are generally two players: a Sender and a Receiver. The Sender has a database DB of records, and the Receiver submits a query  $Q$  to the Sender in order to retrieve a particular record. The main requirement in a SPIR system is privacy for both the Sender and the Receiver. That is, on the one hand the Sender should not

learn any information about the index of the record the Receiver is interested in, and on the other hand, the Receiver should not learn any information about the database, beyond the content of the record defined in the query  $Q$ , and what is already publicly known. In particular, the Receiver should not be able to learn information about more than one record per query. For instance, the Receiver should not be able to learn, through one query, the value of any function on a set of more than one record. SPIR systems have many real-life applications; for instance, consider a scenario where the inventor of a new drug needs information on a number of chemical components that will constitute his final product. This information can be accessed *for a fee* at some central database. This database could be managed, however, by parties with possibly competitive interests, and the inventor fears that his intellectual property (IP) will be compromised. He would like, therefore, that his queries remain concealed from the database manager. The latter, on the other hand, wants to be paid for all information retrieved from his database. It is clear that the SPIR system described above, can be a solution to this set of conflicting requirements.

There are similar applications however, that are closely related to the IP example above, which cannot be solved by a SPIR primitive. Consider for example the following e-health scenario where three types of participants are involved: (1) a patient, (2) a medical database containing the health records of patients, and (3) a doctor querying the medical database on patients' health records. The medical database and the doctor can be thought of as the Sender and Receiver, respectively, in a traditional SPIR setting. The requirements in the e-health application are as follows:

1. **Privacy for the Receiver:** The Receiver (doctor) wants to retrieve records from the medical database, without the Sender (DB) learning the index of those records, and thus the identity of his patient.
2. **Privacy for the Sender:** The Sender (DB) wants to be sure that, for each query, the Receiver (doctor) learns information only on one record (defined in the query) and nothing about the other records.
3. **Privacy for the data subject:** In order to comply with privacy legislation, the Sender wants to be sure that the Receiver has an valid reading authorization from the owner of the targeted record (i.e., the patient). We call the latter, *an Authorizer*. Notice that the Sender should not be able to learn the Authorizer's identity, otherwise the first requirement will be violated.

The example above shows a typical scenario where plain SPIR primitives fall short of protecting the interests of the Sender, the Receiver, and the Authorizer at the same time. The solution we provide in this paper, addresses the interests of all three parties, and solves the problems described above. We call the presented solution: *Accredited SPIR*. In what follows, we sometimes refer to the latter set of requirements, namely privacy for the data-subject, the Sender, and the Receiver, as the *Accredited SPIR problem*.

**SOLUTION HIGHLIGHT.** In the Accredited SPIR setting we have three players: a Sender, a Receiver, and an Authorizer. The Receiver submits a query  $Q$  to the Sender, who replies with a response  $R$ . The Receiver recovers the answer to his query from  $R$ . The main contribution of Accredited SPIR, is to assure the Sender, before processing the query  $Q$ , that the Receiver has obtained an explicit consent from the owner of the record defined in  $Q$ , *without* revealing the identity of this owner (i.e., the Authorizer).

The Accredited SPIR architecture we propose, combines three cryptographic primitives: privacy-preserving digital credentials, homomorphic encryption, and SPIR systems. Privacy-preserving digital credentials [Cha85, CP92, Bra00, CL02, CL04] are cryptographic tokens issued by a certification authority CA to individuals. The CA encodes in each credential a set of attributes about the identity of its recipient. The latter is called a credential holder. A credential holder may later show his credential to a verifier in return for a service or a privilege (e.g., to receive medical treatment). Unlike traditional PKI certificates (e.g., X.509), privacy-preserving digital credentials allow their holders to selectively disclose information about their attributes [Bra00]. In particular, if a credential holder has a set of attributes  $(x_1, \dots, x_n)$ , then he can prove any predicate  $\mathcal{P}$  satisfied by those attributes, without the verifier learning any extra information beyond the status of  $\mathcal{P}(x_1, \dots, x_n)$ .

Assume the Authorizer has a CA-issued identity credential  $Cred$  containing a set of attributes (ID, Age,  $\dots$ ). The idea is to first make the Authorizer and Receiver jointly compute the query  $Q$ , and then have the Authorizer produce a signed proof of knowledge of the secret attributes embedded in  $Cred$ . Along with the latter, the Authorizer proves that the ID attribute embedded in  $Cred$  is the same as the one contained in the query  $Q$ . The Receiver then deposits the signed proof along with the query to the Sender. The Sender first checks the validity of the proof. If accepted it carries on with the SPIR protocol and processes the query, otherwise it rejects.

As mentioned earlier, the signed proof does not reveal any information about the credential holder, and yet guarantees that the content of the query is consistent with the secret identity attribute embedded in the credential. Furthermore, owing to the fact that it is hard for a polytime adversary to forge credentials, or make proofs about credentials he does not own, the Sender can be sure that the Receiver has indeed obtained an explicit consent from the targeted record's owner.

This paper presents three constructions to solve the accredited SPIR problem. The first is based on a modified version of one of Brands DL-based credentials [Bra00, Section 4.5.2], the ElGamal cryptosystem, and a SPIR system proposed by Lipmaa in [Lip05]. The two additional constructions are variants of the first, and use an RSA-based version of Brands credentials [Bra00, Section 4.2.2], in combination with the ElGamal, and the Okamoto-Uchiyama [OU98] cryptosystems. In the following, we describe previous results and related work available in the literature.

## 2 Related Work

Much research has gone into the problem of managing personal data in accordance with a user-defined privacy policy. In [GMM06], for instance, Golle *et al.* propose a mechanism by which data collectors can be caught and penalized if they violate an agreed-upon policy, and disclose sensitive data about a data-subject. The main idea there is that a data-collector would place a *bounty*, which it must forfeit if a privacy violation is uncovered. The bounty could be explicit in the form of a bond, or implicit in the form of penalties imposed if privacy is violated. This technique however is geared towards violation detection after the fact, and assumes the existence of active *bounty hunters* who seek to induce dishonest data collectors into committing unlawful disclosures.

Another related approach is that of policy-based encryption by Bagga *et al.* [BM05, BM06]. Policy-based encryption allows a user to encrypt a message with respect to a credential-based policy, formalized as a monotone boolean expression. The encryption is such that only a user having access to a qualified set of credentials, complying with the policy, is able to successfully decrypt the message. The context in [BM05, BM06], however, is different from the one in this paper, since the goal there is to allow the user to send a secret message to a designated set of players defined by a policy. In our context, the user's data is already stored in a database, and the goal is to allow user-authorized parties to retrieve the user's data, without the database manager learning which data has been retrieved or the identity of the data subject. It is also not clear how revocability can be implemented in the context of [BM05, BM06].

In [SWP00], Song *et al.* present a scheme for the problem of searching keywords on encrypted data. The setting there consists of a user, and a server storing encrypted data owned by the user. The server can process search queries on the user's stored ciphertext, only if given proper authorization from the user. The proposed scheme also supports hidden user queries, where the server conducts the search without learning anything about the content of the query. Although somewhat related to our context, it is not clear how the work in [SWP00] can be applied to the problem we describe in this paper, since delegating querying capabilities to a third party, may require the user to reveal his encryption key, and thus share all of his past and future secrets. Besides, it is not clear how the identity of the data-owner can be hidden from the server, or how to impose (e.g., time or usage) restrictions on search capabilities delegated to a third party.

Finally, in [AIR01] Aiello *et al.* consider a scenario where a database contains a set of priced data items, and users privately retrieve data from it. The proposed protocol is called priced oblivious transfer, and allows a user  $U$ , who made an initial deposit, to buy different data items, without the database manager learning which items  $U$  is buying, as long as  $U$ 's balance contains sufficient funds. We believe the construction in [AIR01] is the first to impose additional requirements on oblivious transfer protocols. While interesting in their own right, the added requirements do not address the identity of the data-subjects, and hence do not seem to help in solving the problem we consider in this paper.

### 3 Summary of Contribution and Paper Organization

We propose a solution to the accredited SPIR problem, and three constructions to implement it. The solution we present allows a user to issue authorizations to a Receiver to privately retrieve his records, without the database manager learning anything about the retrieved data, or the data subject's identity. The authorizations contain computationally non-modifiable, unforgeable, user-defined policies and limitations, governing their use. The authorizations can also be anonymously revoked by their issuer if needed.

To the best of our knowledge, this work is the first to give a solution to the accredited SPIR problem, and to address the more general issue of enforcing user-defined privacy policies, by combining SPIR protocols and privacy-preserving digital credentials.

The remainder of this paper is organized as follows. In Section 4, we describe the main building blocks used in the first construction and throughout the paper. In Section 5, we present a DL-based accredited SPIR construction. In Sections 6 and 7, we discuss the security, privacy features, and performance of the first construction. In Section 8.2, we present a second construction based on a RSA version of Brands credentials. In Section 8.4, we give a third variant based on the Okamoto-Uchiyama cryptosystem. We conclude in Section 9.

## 4 Building Blocks for the DL-Based Construction

### 4.1 Brands-CP Credentials

In [Bra00], Brands proposes various credential systems based on the hardness of the discrete logarithm problem in groups of prime order, and the RSA problem in groups of composite order (RSA groups). Brands has also proposed other variants of the above systems, based on DSA and the Chaum-Pedersen signatures [Bra00, Section 4.5.2]. For the purpose of our first construction, we will use the latter variant, and will refer to it as the Brands-CP system. The security of the Brands-CP system is based on the hardness of the discrete logarithm problem in groups of prime order. The Brands-CP system allows a certification authority CA to issue to a user U a set of credentials encoding attributes about U's identity. The credential itself consists of (1) a public key  $h$  embedding the user's attributes and (2) a special CA-supplied digital signature on it, denoted  $\sigma_{CA}(h)$ . At the end of the issuing protocol, the credential that user U has obtained is perfectly hidden from the CA, and perfectly indistinguishable from any other credential the CA has previously issued. Later, user U can show his credential to individuals and organizations in return for a service. Showing a credential does not necessarily require the revealing of the attributes encoded in it. A credential holder can selectively and verifiably disclose any information he wishes about his attributes, which may include revealing the actual values of the attributes, or just proving a predicate about them. In [Bra00, Section 3.6], Brands shows how to prove a class of linear predicates about the attributes. At a later stage, and depending on the application, the verifying individual or organization may want to deposit

the credential showing transcript to the certification authority. This deposit can be thought of as a cheque deposit in the context of e-banking or as a ballot submission in the context of e-voting. The deposited transcript is unlinkable to the instance of the issuing protocol that generated the credential. For the sake of completeness, we give in the following a brief description of the issuing, and showing protocols of the Brands-CP system, as well as an overview of the parameters and setting.

*System setting.* On input the security parameter  $\kappa$ , the CA chooses  $\kappa$ -sized primes  $p$  and  $q$  such that  $2q|p - 1$ . Let  $G_q$  be the unique subgroup of  $\mathbb{Z}_p^*$  of order  $q$ , and let  $g_0$  be one of its generators. The CA also chooses  $H : \{0, 1\}^* \rightarrow \mathbb{Z}_q^*$ , a public collision-resistant hash function. In the setup phase, the certification authority randomly chooses  $y_1, y_2, \dots, y_\ell$  and  $x_0 \in_R \mathbb{Z}_q^*$ , and computes  $(g_1, g_2, \dots, g_\ell, h_0) := (g_0^{y_1}, g_0^{y_2}, \dots, g_0^{y_\ell}, g_0^{x_0}) \pmod p$ . The parameters  $(g_1, g_2, \dots, g_\ell, h_0)$  are then made public along with  $g_0, q$ , and  $G_q$ .

*Credential issuing.* To obtain a credential, a user first convinces the certification authority that he fulfills a set of application-specific requirements necessary to receive that credential. The certification authority then encodes a set of  $\ell$  of user attributes in the credential. Let  $x_1, \dots, x_\ell$  denote the attributes to be encoded. The credential's public key is then computed as  $h := (g_1^{x_1} \dots g_\ell^{x_\ell} h_0)^\alpha$ , where  $\alpha$  is a secret blinding factor randomly chosen in  $\mathbb{Z}_q^*$  by the user. The certification authority's signature on the credential is a triplet  $(c'_0, r'_0, z') \in \mathbb{Z}_q^2 \times G_q$ , satisfying the relation  $c'_0 = H(h, z', g_0^{r'_0} h_0^{-c'_0}, h^{r'_0} z'^{-c'_0})$ . At the end of the issuing protocol, the certification authority knows neither  $h$  nor the signature  $(c'_0, r'_0, z')$ .

*Credential showing.* In order to have access to a service, user U can show his credential without the verifying party being able (1) to learn information about the encoded attributes beyond what U willingly discloses, or (2) to link the credential to the user's identity even if it colludes with the certification authority. In practice, to show his credential to a verifying party, user U reveals (1) the credential's public key  $h$  along with a signature  $\sigma_{CA(h)} := (z', c'_0, r'_0)$ , and (2) a signed proof of knowledge of a representation of  $h$ , with respect to basis  $(g_1, g_2, \dots, g_\ell, h_0)$ . This signature is performed on a verifier-chosen challenge  $m$ . The verifier checks the validity of the credential by verifying if the relation  $c'_0 \stackrel{?}{=} H(h, z', g_0^{r'_0} h_0^{-c'_0}, h^{r'_0} z'^{-c'_0})$  holds. If the credential is valid, the verifier moves on to check the validity of the signed proof of knowledge.

More details on the Brands-CP credential issuing and showing protocols, as well as their security and privacy properties, can be found in [Lay07].

## 4.2 ElGamal Homomorphic Encryption

Our DL-based accredited SPIR construction relies on the ElGamal encryption scheme because of its homomorphic properties and because it fits well the setting of the Brands-CP credentials. In the following we recall the settings of the ElGamal cryptosystem.

*Settings.* Let  $p, q,$  and  $G_q,$  be the public parameters chosen by the CA in the setup of the Brands-CP credential system. User  $U$  randomly chooses  $g_{\text{ElG}}$ , a generator of  $G_q,$  and  $x_u \in_R \mathbb{Z}_q^*,$  and computes  $y_{\text{ElG}} := g^{x_u} \bmod p.$  User  $U$  then publishes his ElGamal public key  $(G_q, g_{\text{ElG}}, y_{\text{ElG}}),$  and keeps his private key  $x_u$  secret. A message  $m \in G_q,$  can be encrypted by choosing a random  $r \in_R \mathbb{Z}_q^*,$  and computing  $c = (g_{\text{ElG}}^r, y_{\text{ElG}}^r m) = (c_1, c_2).$  Using  $U$ 's private key, the plaintext can be recovered as  $m = c_2/c_1^{x_u}.$  Given a constant  $\alpha,$  and encryptions of  $m$  and  $m',$  it is easy to compute randomized encryptions of  $m \times m'$  and  $m^\alpha.$

### 4.3 AIR-Based Lipmaa $\binom{1}{n}$ -OT

In [Lip05], Lipmaa proposes a SPIR scheme based on ideas from a construction by Aiello et al. [AIR01]. Lipmaa's SPIR scheme is computationally private for the Receiver and perfectly private for the Sender. Its security relies on the hardness of the *decisional Diffie-Hellman* and the *decisional composite residuosity* problems [Lip05]. The SPIR scheme in [Lip05] has a log-squared communication complexity (in the size of the database).

**Main idea.** Let  $DB$  denote the Sender's private database, and let  $s$  be the index of the record the Receiver is interested in. The receiver computes  $c := E_{pk_{\text{hom}}}(s),$  a homomorphic encryption of  $s,$  and sends it to the Sender. Using the homomorphic properties of the encryption, the Sender computes for each record  $DB[j]$  in the database,  $DB'[j] := E_{pk_{\text{hom}}}(\delta_j(s - j) + DB[j]),$  where  $\delta_j$  is a random blinding factor chosen by the Sender. The encrypted records  $DB'[j]$  are then sent to the Receiver, who will be able to retrieve something meaningful only from  $DB'[s] := E_{pk_{\text{hom}}}(DB[s]);$  everything else will decrypt to randomness. The construction in [Lip05] follows a similar methodology to the above, except that the Sender uses an extra loop of superposed encryptions that leads to a randomized ciphertext of  $DB'[s].$  Only the latter is sent back to the Receiver. This is done as follows. The database  $(DB[1], \dots, DB[n])$  is arranged in an  $\alpha$ -dimensional  $\lambda_1 \times \dots \times \lambda_\alpha$  hyper-rectangle for some pre-defined positive integers  $\lambda_j,$  such that  $n = \prod_{j=1}^\alpha \lambda_j.$  Each record  $DB[i]$  is indexed by a tuple  $(i_1, \dots, i_\alpha)$  on this hyper-rectangle, where  $i_j \in \mathbb{Z}_{\lambda_j}.$  To retrieve a particular record  $(s_1, \dots, s_\alpha),$  the Receiver submits to the Sender a homomorphic encryption  $\beta_{jt} := HomEnc_{pk}(b_{jt}),$  for  $1 \leq j \leq \alpha, 0 \leq t < \lambda_j,$  where  $b_{jt} = 1$  if  $t = s_j,$  and  $b_{jt} = 0$  otherwise. The Sender exploits the homomorphic properties of the encryption scheme  $HomEnc_{pk}(\cdot)$  to create a new  $(\alpha - 1)$ -dimensional database  $DB_1,$  such that  $\forall (i_2, \dots, i_\alpha) \in \mathbb{Z}_{\lambda_2} \times \dots \times \mathbb{Z}_{\lambda_\alpha},$   $DB_1(i_2, \dots, i_\alpha)$  is equal to an encryption of  $DB_0(s_1, i_2, \dots, i_\alpha),$  where  $DB_0$  is the the Sender's original database  $DB.$  The same procedure is repeated, and at the  $j^{\text{th}}$  iteration, an  $(\alpha - j)$ -dimensional database  $DB_j$  is obtained by the Sender, such that  $DB_j(i_{j+1}, \dots, i_\alpha)$  is equal to a  $j$ -times encryption of  $DB_0(s_1, \dots, s_{j-1}, i_j, \dots, i_\alpha).$  After  $\alpha$  iterations, the Sender obtains  $DB_\alpha,$  an  $\alpha$ -times encryption of  $DB_0(s_1, \dots, s_\alpha).$  The Sender returns  $DB_\alpha$  to the Receiver, who needs to decrypt it  $\alpha$  times to recover  $DB(s).$  In [Lip05], Lipmaa uses a

special *Length Flexible Additively Homomorphic* encryption scheme to implement  $HomEnc_{pk}(\cdot)$ . Notice that in the hyper-rectangle construction above, the Receiver can cheat by maliciously sending  $\beta_{jt} := HomEnc_{pk}(1)$ , for  $t \neq s_j$ . To stop such attacks, the Sender performs the repeated encryptions above, on  $DB_0 = DB'$  rather than on the Sender's original DB. Let  $s'$  be the index corresponding to the  $\beta_{jt}$ 's. At the end of the protocol, the Receiver obtains  $DB_\alpha$ , which he decrypts  $\alpha$  times to recover  $DB'[s'] = E_{pk_{hom}}(\delta_{s'}(s - s') + DB[s'])$ . Next, the Receiver decrypts  $DB'[s']$  once again, and recovers something meaningful ( $DB[s]$ ) only if  $s' = s$ . More details can be found in [Lay07].

**Remark.** Both [Lip05] and [AIR01] propose the ElGamal cryptosystem to implement  $E_{pk_{hom}}$ . It is worth noting however, that using plain ElGamal, it is not possible to compute  $E_{pk_{hom}}(\delta_j(s - j) + DB[j])$  given  $E_{pk_{hom}}(s)$ , since ElGamal is *only multiplicatively* homomorphic. We fix this problem in the next section.

### 5 Accredited SPIR Based on the DL Problem

The DL-based accredited SPIR scheme we propose, is achieved by combining the three building blocks above, modulo few adaptations. We first give a high-level overview of the construction, before getting into the details. We assume the public parameters of the three building blocks are already known to all parties. Let  $ID_A$  be an attribute, that uniquely identifies the Authorizer (e.g., an SSN). This  $ID_A$  will determine the index by which the Receiver will query the Sender's database. Let us first assume that the Authorizer possesses a Brands-CP credential of the form  $(h, \sigma_{CA}(h))$ , where  $h = (g_1^{ID_A} g_2^{x_2} \dots g_\ell^{x_\ell} h_0)^\alpha$ . The Authorizer computes  $c := E_{pk_{hom}}((g_{db})^{ID_A}) := E_{pk_{EIG}}((g_{db})^{ID_A}) := (c_1, c_2)$ , where  $pk_{EIG}$  is the Receiver's ElGamal public key, and  $g_{db}$  is a public generator of  $G_q$  chosen by the Sender. Next, the Authorizer produces a signed proof of knowledge asserting that the logarithm, to the base  $g_{db}$ , of the plaintext encoded in  $c$ , is the same as the first attribute embedded in credential  $h$ . We call this last assertion an ID-consistency proof.

Notice that this latter proof cannot be done in a straightforward way using the original Brands-CP credentials, because  $h$  has the form  $h := g_1^{\beta_1} \dots g_\ell^{\beta_\ell} h_0^\alpha$ , where  $\beta_i = \alpha x_i$  for some random blinding factor  $\alpha$ . Establishing ID-consistency in this case requires proving a non-linear predicate on secret exponents  $(\alpha, \beta_1, ID_A)$ , defined by  $\mathcal{P} \equiv \beta_1 = \alpha \times ID_A$ , which cannot be done efficiently. To fix this problem we propose a modified version of the Brands-CP credentials, with exactly the same security and privacy properties. In the modified version, the credential's public key  $h$  is computed as  $h := (g_1^{x_1} \dots g_{\ell-1}^{x_{\ell-1}} g_\ell^\alpha h_0)$ , where  $x_1, \dots, x_{\ell-1}$  are identity attributes, and  $\alpha$  is a secret random blinding factor chosen by the credential recipient. This modification is of general interest, and can be used in other contexts as well. For lack of space, a summary of the modified credential system is outlined in [Lay07].

In what follows, we assume the Authorizer possesses a credential of the new type. Let us denote the public key of the Authorizer's new credential by  $h := (g_1^{ID_A} \dots g_{\ell-1}^{x_{\ell-1}} g_\ell^{\alpha_1} h_0)$ . To prove ID-consistency between  $h$  and the SPIR query

$(c_1, c_2) := E_{\text{pk}_{\text{ElG}}}((g_{db})^{\text{ID}_{\mathcal{A}}})$ , it suffices for the Authorizer to produce a signed proof of knowledge of a DL-representation of  $(h/c_2 \bmod p)$  with respect to basis  $((g_1 g_{db}^{-1}), g_2, \dots, g_\ell, h_0, g_{\text{ElG}}, y_{\text{ElG}})$ . We denote the latter signed proof by  $\text{SPK}\{(\varepsilon_1, \dots, \varepsilon_\ell, \mu, \nu) : h = g_1^{\varepsilon_1} \dots g_\ell^{\varepsilon_\ell} h_0 \wedge c_2 = y_{\text{ElG}}^\mu g_{db}^\nu \wedge \varepsilon_1 = \nu\}(m)$ . The message  $m$  to be signed, can be a concatenation of several fields, including a fresh nonce. In addition, it may contain the identity of the Receiver, which will allow the Authorizer to exclusively tie the authorization to the Receiver, and discourage him from sharing it with a third party. The Authorizer may also include an expiry date in  $m$  to make sure his authorization remains valid only for the appropriate amount of time. More generally, the Authorizer may encode in  $m$  any application-specific policy he wants the Receiver to follow.

**Running the SPIR.** Now let us assume the signed proof above was accepted. The next step, would be for the Receiver to compute the query messages as indicated in the OT scheme of Section 4.3. First let  $(\text{ID}_{\mathcal{A}}^{(1)}, \dots, \text{ID}_{\mathcal{A}}^{(\alpha)})$  be the representation of the Authorizer's  $\text{ID}_{\mathcal{A}}$  in the  $\alpha$ -dimensional hyper-rectangle  $\lambda_1 \times \dots \times \lambda_\alpha$  used by the Sender's database. The Receiver then computes, for  $1 \leq j \leq \alpha, 0 \leq t < \lambda_j$ , the homomorphic encryptions  $\beta_{jt} := \text{HomEnc}_{\text{pk}}(b_{jt})$ , where  $b_{jt} = 1$  if  $t = \text{ID}_{\mathcal{A}}^{(j)}$ , and  $b_{jt} = 0$  otherwise. Next, the Receiver submits to Sender: (1) the credential  $(h, \sigma_{\text{CA}}(h))$ , (2) the first part of the query  $(c_1, c_2) := E_{\text{pk}_{\text{ElG}}}((g_{db})^{\text{ID}_{\mathcal{A}}})$ , (3) an ID-consistency proof  $\text{SPK}\{(\varepsilon_1, \dots, \varepsilon_\ell, \mu, \nu) : h = g_1^{\varepsilon_1} \dots g_\ell^{\varepsilon_\ell} h_0 \wedge c_2 = y_{\text{ElG}}^\mu g_{db}^\nu \wedge \varepsilon_1 = \nu\}(m)$ , and (4) the second part of the query consisting of the  $\beta_{jt}$ 's for  $1 \leq j \leq \alpha, 0 \leq t < \lambda_j$ .

Note that there is no need for the Receiver to prove consistency between the  $\beta_{jt}$ 's and  $(c_1, c_2)$ . As we will show later, any attempt by the Receiver to incorrectly compute the  $\beta_{jt}$ 's, will prevent him from learning anything meaningful at the end of the SPIR protocol.

Once the ID-consistency check succeeds, the Sender starts processing the Receiver's query as explained in the following. But first, we make few practical assumptions. We assume that  $n$ , the size of the Sender's database, is bounded above by  $q$ , the order of  $G_q$ . In practice,  $q$  is chosen to be at least 160-bit long, which means the Sender's database could have up to  $2^{160}$  different records. Although we think this should be sufficient in practice, the size of  $q$  can always be increased if needed. Moreover, we assume that each record  $\text{DB}[i]$  contains a field for storing  $(g_{db}^{-i} \bmod p)$ , in addition to a large field containing application-specific data (e.g., health, financial data).

Now the query is processed as follows. Using the first part of the query, and the *multiplicatively* homomorphic properties of ElGamal, the Sender computes for  $j \in [1, n]$ ,  $\text{DB}_0[j] = E_{\text{pk}_{\text{ElG}}}((g_{db})^{\delta_j (\text{ID}_{\mathcal{A}} - j)} \times \text{DB}[j]) = ((E_{\text{pk}_{\text{ElG}}}(g_{db}^{\text{ID}_{\mathcal{A}}}) \times g_{db}^{-j})^{\delta_j} \times \text{DB}[j])$ , where  $\delta_j$  is a random blinding factor chosen by the Sender. Note that  $g_{db}^{-j}$  has already been precomputed and stored in with  $\text{DB}$ 's  $j^{\text{th}}$  record. The Sender then proceeds with computing  $\text{DB}_\alpha$  by repeated encryptions of the records of  $\text{DB}_0$  as indicated in the OT scheme of Section 4.3. Upon receiving  $\text{DB}_\alpha$ , the Receiver recovers  $\text{DB}_0[\text{ID}_{\mathcal{A}}]$  by repeated decryption using his secret homomorphic key  $sk$ . Next the sender obtains the desired record  $\text{DB}[\text{ID}_{\mathcal{A}}]$  by decrypting

$DB_0[ID_{\mathcal{A}}]$  using his ElGamal private key. A figure summarizing the protocol can be found in [Lay07].

## 6 Security and Privacy Properties

**Theorem 1.** *Assuming the DL problem is hard, the ElGamal and  $HomEnc(\cdot)$  cryptosystems IND-CPA secure, and that the none of the three parties colludes with the other, the protocol of Section 5 solves the Accredited SPIR problem in the random oracle model, and provides computational privacy for both the Authorizer and Receiver, and perfect privacy for the Sender.*

*Proof Sketch.*<sup>1</sup> *Correctness.* easy to check and relies on the homomorphic properties of the ElGamal and the  $HomEnc(\cdot)$  cryptosystems.

*Soundness.* Assume the Receiver does not follow the protocol, and maliciously uses an index  $s \neq ID_{\mathcal{A}}$  in the second part of the query. This will lead  $DB_{\alpha}$  computed by the Sender to be an  $\alpha$ -time encryption under  $HomEnc_{pk}(\cdot)$  of  $DB_0[s] = E_{pk_{EG}}((g_{db})^{\delta_s (ID_{\mathcal{A}} - s)} \times DB[s])$ , where  $\delta_s$  is a secret blinding factor.  $DB_0[s]$  decrypts to random, which results in perfect privacy for the Sender.

*Privacy for the Receiver.* The Sender's view consists of a credential, an ElGamal encryption, a signed proof of knowledge, and a  $HomEnc(\cdot)$  encryption of an  $\alpha$ -dimensional coordinate. Because the proof of knowledge is honest-verifier zero-knowledge for any distribution of the attributes (cf. [Bra00, Prop. 3.3.4]), seeing the signed proof of knowledge, together with the credential, does not leak any information to the Sender about the content of the Receiver's query. The Sender could only hope to extract information from the ElGamal ciphertext  $c$  and the homomorphic encryptions  $\beta_{jt}$ 's. But this should not be possible, because it implies breaking the security of ElGamal, or the  $HomEnc(\cdot)$  cryptosystems, which contradicts our assumption.

*Privacy for the Authorizer.* Assuming the Receiver and the Sender do not collude, there are two ways to violate the Authorizer's privacy. Either by learning information from the Receiver's query, or by forging a signed proof of knowledge on behalf of the Authorizer. The former attack is computationally impossible and follows from the Receiver's privacy. To achieve the second however, one could either (1) forge a new credential from scratch with the Authorizer's identity embedded in it, (2) forge a signed proof on a legitimately-obtained credential that was not issued to the Authorizer, and thus does not initially contain his identity, or (3) forge a signed proof on behalf of the Authorizer on a credential owned by the Authorizer. Attack (1) can be ruled out based on the computational unforgeability of the Brands-CP credentials in the random oracle model (cf. [Bra00, Prop. 4.3.7]). Similarly, attack (2) is impossible because of the *non-modifiability* property of Brands-CP credentials [Bra00, Prop. 3.3.8].

---

<sup>1</sup> Complete details will be given in the full version of the paper.

The non-modifiability feature states that, assuming the DL problem is hard, it is infeasible, in the random oracle model, to construct a signed proof for a formula that does not in fact apply to the prover's representation. Finally, attack (3) can be ruled out as well, owing to the unforgeability of signed proofs on Brands-CP credentials [Bra00, Prop. 3.3.6].

## 6.1 Additional Privacy for the Authorizer

*User-centricity and policy enforcement.* When issuing an authorization, the credential holder could specify in the message of the signed proof of knowledge, a set of rules that he wishes the authorization recipient to comply with. For instance, he could specify an expiry date, an upper bound on the number of times the authorization is used, or any other usage policy. The Sender is supposed to refuse processing queries from a Receiver who does not satisfy the usage policy specified in the signature.

*Revocability.* The Authorizer may decide to revoke a previously issued authorization. This can be done anonymously as follows. The Authorizer first needs to prove knowledge, over an physically anonymized channel (e.g., a MixNet [Cha81]), of a representation of the credential used in the authorization to be revoked. Once the proof of knowledge is accepted, the Authorizer requires the Sender (DB manager) to add the credential in question to a black list of revoked authorizations. Later, it is easy for the Sender to check whether the credential contained in a submitted query is on the black list or not. This can be done efficiently using hash tables for instance. Note that an authorization can be revoked only by its issuer, since an polytime adversary cannot find a discrete-log representation with non-negligible probability.

*Authenticated personal information retrieval.* In the special case, where the Authorizer and Receiver are the same entity, the construction we propose provides the data-subject with a mechanism to retrieve *his own* personal data anonymously. Our construction also ensures that the stored data can be retrieved only by its owner. The channel between the Receiver and Sender in this case has to be physically anonymized.

## 7 Performance Analysis

The accredited SPIR construction of Section 5 does not lead to a significant increase in computation and communication complexity, compared to the underlying SPIR scheme [Lip05]. If we assume the Authorizer has a credential with  $(\ell - 1)$  attributes, then the added computation complexity is as follows. The Authorizer needs to make  $(\ell + 6)$  offline exponentiations (all precomputable), while the Receiver and Sender, both need to make  $(\ell + 8)$  online exponentiations. This is negligible compared to the complexity of the underlying SPIR scheme which is linear in  $n$ , the size of the database. In practice,  $\ell$  is in the range of 20, while  $n \approx 2^{160}$ . In terms of communication complexity, both the

Authorizer and Receiver need to send  $(\ell + 8) \log(n) + 5$  extra bits to the Receiver and Sender respectively. Again this does not change the overall  $\mathcal{O}(\log^2(n))$  asymptotic communication complexity of the underlying SPIR scheme [Lip05].

## 8 Accredited SPIR Based on RSA

The constructions we present in this section are based on a RSA-version of Brands’ credentials [Bra00, Section 4.2.2]. For the sake of completeness, we briefly introduce them in the following.

### 8.1 Brands-RSA Credentials

*Settings.* On input the security parameter  $\kappa$ , the credential issuer chooses: (1)  $\kappa$ -sized primes  $P$  and  $Q$ , and computes  $N := PQ$ , (2) a prime  $v$  smaller than  $N$ , and co-prime to  $\phi(N)$ , (3) random elements  $(g_1, \dots, g_\ell) \in_R (\mathbb{Z}_N^*)^\ell$ , and (4) a one-way hash function  $\mathcal{H}(\cdot) = \mathcal{H}_{N,v}(\cdot) : \{0, 1\}^* \rightarrow \mathbb{Z}_s$ , for some  $s$  superpolynomial in  $\kappa$ .

The credential issuer makes the parameters  $N, v, (g_1, \dots, g_\ell), \mathcal{H}$  public, and keeps  $P$  and  $Q$  secret. In addition, the issuer chooses  $x_0 \in_R \mathbb{Z}_v^*$ , such that given  $h_0 := x_0^v \pmod N$ , computing the  $v^{\text{th}}$  root of  $h_0$  is hard. The issuer then publishes  $h_0$  and keeps  $x_0$  secret.

*Credential issuing.* Assume after making the necessary identity checks, the certification authority accepts to issue a credential to the user. Let  $(x_1, \dots, x_\ell) \in (\mathbb{Z}_v^*)^\ell$  be the attributes the CA wants to encode in the credential, and let  $h := g_1^{x_1} \dots g_\ell^{x_\ell} \pmod N$ . The  $x_i$ ’s are known to both the user and the CA. The user then chooses a random blinding factor  $\alpha_1 \in \mathbb{Z}_N^*$  and computes the credential’s public key  $h' := h\alpha_1^v$ . The certification authority’s digital signature on the credential is a pair  $(c'_0, r'_0) \in \mathbb{Z}_s \times \mathbb{Z}_N^*$ , satisfying the relation  $c'_0 = \mathcal{H}(h', r_0^v (h_0 h')^{-c'_0})$ . At the end of the issuing protocol, the certification authority knows neither  $h'$  nor the signature  $(c'_0, r'_0)$ .

*Credential showing.* Similar to the Brands-CP system, a user can show his credential to a verifying party, by first revealing the credential’s public key  $h'$  and CA-signature  $(c'_0, r'_0)$ . The verifier checks if the validity relation  $c'_0 \stackrel{?}{=} \mathcal{H}(h', r_0^v (h_0 h')^{-c'_0})$  holds. Once the validity check succeeds, the user produces a signed proof of knowledge of a RSA representation  $(x_1, \dots, x_\ell, \alpha_1)$  of  $h'$  with respect to basis the  $(g_1, \dots, g_\ell, v)$ . The signed proof can also be computed with respect to a predicate  $\mathcal{P}$  on exponents  $(x_1, \dots, x_\ell)$ , agreed-upon by the user and the verifier at the time of the showing.

More details on the Brands-RSA credential issuing and showing protocols can be found in [Lay07].

### 8.2 Combining ElGamal Encryption with Brands RSA-Based Credentials

We assume the Authorizer possesses a Brands-RSA credential  $h'$ , and certificate  $\sigma_{\text{CA}}(h') := (c'_0, r'_0)$ , with  $h'$  of the form  $h' = (g_1^{\text{ID}_A} g_2^{x_2} \dots g_\ell^{x_\ell} \alpha^v) \pmod N$ . Recall

that  $ID_{\mathcal{A}}, x_2 \cdots x_\ell$  are elements of  $\mathbb{Z}_v$ , and that  $ID_{\mathcal{A}}$  represents a record index in the Sender’s DB. To accommodate all possible DB indexes  $ID_{\mathcal{J}}$  (ranging over  $[1, n] = [1, q]$ ), the prime  $v$  is chosen to be greater than  $q$ , the order of  $G_q$  in the ElGamal setting. Also, for reasons that will become clear shortly, the prime  $v$  is chosen to be coprime to  $p - 1$ , where  $p$  denotes the public prime parameter in the ElGamal setting. Finally, the RSA factors  $P$  and  $Q$  are chosen to be greater than  $p$ .

As in the first construction based on Brands-CP credentials, the Authorizer (data subject) and Receiver use the ElGamal cryptosystem to compute the SPIR query. Assuming the same setting for the Sender and Receiver as in Section 5, the Authorizer computes  $(c_1, c_2) := E_{pk_{\text{ElG}}}((g_{db})^{ID_{\mathcal{A}}})$ . To prove ID-consistency between  $h'$  and the SPIR query, it suffices for the Authorizer to produce a signed proof of knowledge of a RSA-representation of  $(h'/c_2 \pmod{Np})$  with respect to basis  $((g_1 g_{db}^{-1}), g_2, \dots, g_\ell, y_{\text{ElG}}^{-1}, v)$ .

*Putting the pieces together.* As in the Brands-CP case, the Authorizer proves ID-consistency between his credential and the query by sending to the Receiver a signed proof of knowledge of the form  $\text{SPK}\{(\varepsilon_1, \dots, \varepsilon_{\ell+1}, \mu) : h'/c_2 = (g_1 g_{db}^{-1})^{\varepsilon_1} g_2^{\varepsilon_2} \cdots g_\ell^{\varepsilon_\ell} (y_{\text{ElG}}^{-1})^{\varepsilon_{\ell+1}} \mu^v \pmod{Np}\}(m)$ . The Authorizer can use the message  $m$  to encode any usage policy he wants the Receiver to follow.

After accepting the signed proof, the Receiver proceeds with the SPIR protocol of Section 5 without any further changes.

### 8.3 Security and Privacy Properties

**Theorem 2.** *Assuming the (multi-prime) RSA problem is hard, the ElGamal and HomEnc(.) cryptosystems IND-CPA secure, and that none of the three parties colludes with the other, the protocol of Section 8.2 solves the Accredited SPIR problem in the random oracle model, and provides computational privacy for both the Authorizer and Receiver, and perfect privacy for the Sender.*

Complete details of the proof will be given in the full version of the paper. The properties of *user-centricity*, *revocability*, and *authenticated PIR* described in Section 6.1, do apply for the new scheme as well.

### 8.4 Variant Based on the Okamoto-Uchiyama Cryptosystem

The construction of Section 8.2, can be modified by using the Okamoto-Uchiyama cryptosystem [OU98] instead of ElGamal. The Okamoto-Uchiyama cryptosystem is a probabilistic public key cryptosystem whose security is equivalent to the problem of factoring moduli of the form  $n = p^2q$ , for  $p$  and  $q$  prime. The Okamoto-Uchiyama cryptosystem is additively homomorphic.

*Setting of the Okamoto-Uchiyama cryptosystem.* Given security parameter  $\kappa$ , choose  $\kappa$ -sized primes  $p$  and  $q$ , and let  $n = p^2q$ . The choice of  $p$  and  $q$  should be such that  $\gcd(p, q - 1) = \gcd(q, p - 1) = 1$ . Choose random  $g \in \mathbb{Z}_n^*$ , such that  $g_p = g^{p-1} \pmod{p^2}$  has order  $p$ . Let  $h = g^n \pmod{n}$ . The tuple  $(n, g, h, \kappa)$

is published as the public key, while  $(p, q)$  are kept secret. To encrypt a message  $0 < m < 2^{\kappa-1}$ , select random  $r \in \mathbb{Z}_n$ , and compute  $E_{\text{pk}_{\text{OU}}}(m, r) := g^m h^r \pmod n$ . The decryption function uses a special "logarithmic function" [OU98].

*Putting the pieces together.* The Authorizer uses the Receiver's public key (in this case the Okamoto-Uchiyama public key) to produce the SPIR query and prove ID-consistency between the latter and an Authorizer's credential. Let  $c := E_{\text{pk}_{\text{OU}}}(\text{ID}_A, r) = g^{\text{ID}_A} h^r \pmod n$  be a randomized encryption of the Authorizer's ID. Moreover, let  $(h', \sigma_{\text{CA}}(h'))$  be the Authorizer's Brands-RSA credential, with  $h' = (g_1^{\text{ID}_A} g_2^{x_2} \cdots g_\ell^{x_\ell} \alpha^v) \pmod N$ . The Authorizer computes  $h'/c = ((g_1 g^{-1})^{\text{ID}_A} g_2^{x_2} \cdots g_\ell^{x_\ell} h^{-r} \alpha^v) \pmod Nn$ , and produces a signed proof of knowledge of a RSA representation of  $h'/c$  with respect to basis  $((g_1 g^{-1}), g_2, \dots, g_\ell, h^{-1}, v)$ . With very high probability, parameters  $g$  and  $h$  are coprime with  $N$ , otherwise they can be used to factor  $N$  and break the security of the Brands-RSA credential system. Therefore, with very high probability,  $g^{-1}$  and  $h^{-1}$  exist modulo  $Nn$ . Similarly, with very high probability  $\gcd(v, \phi(Nn)) = 1$ , otherwise  $v$  can be used to factor  $n$  and break the Okamoto-Uchiyama cryptosystem. Therefore, the RSA representation above is well defined. Once the signed proof is accepted, the Receiver deposits the SPIR query together with the signed proof to the Sender. The Sender in turn checks the validity of the credential, and signed proof, and proceeds with the remaining steps of the original SPIR scheme of Section 4.3.

## 9 Conclusion

The paper describes a new access control scheme, where access policies are defined by the data subjects. More specifically, the proposed scheme allows database managers to be convinced that each of their stored data is being retrieved according to the policies of the data subjects, without learning which data has been retrieved or the identity of its owner. We present three constructions based on the discrete logarithm and RSA problems. The constructions we propose rely on anonymous authorizations, and combine SPIR systems and privacy-preserving digital credentials. The authorizations contain non-modifiable, unforgeable, user-defined policies governing their use. Moreover, authorizations can be anonymously revoked by their issuers. The proposed solutions yield only a negligible increase in complexity compared to that of the underlying SPIR scheme. This work can be extended in a number of ways. For example it would be interesting to add a mechanism to support the "authorized and anonymous editing of records". One could also try to improve efficiency, and propose additional constructions based on other building blocks and assumptions.

**Acknowledgement.** The author is grateful to Stefan Brands for interesting discussions and for pointing out related literature. This work is supported by the IWT SBO ADAPID project (Advanced Applications for e-ID cards in Flanders). The author was previously funded in part by the Universitary Mission of Tunisia in North America.

## References

- [AIR01] Aiello, W., Ishai, Y., Reingold, O.: Priced oblivious transfer: How to sell digital goods. In: Pfitzmann, B. (ed.) EUROCRYPT 2001. LNCS, vol. 2045, pp. 119–135. Springer, Heidelberg (2001)
- [BM05] Bagga, W., Molva, R.: Policy-based cryptography and applications. In: Patrick, A.S., Yung, M. (eds.) FC 2005. LNCS, vol. 3570, pp. 72–87. Springer, Heidelberg (2005)
- [BM06] Bagga, W., Molva, R.: Collusion-free policy-based encryption. In: Katsikas, S.K., Lopez, J., Backes, M., Gritzalis, S., Preneel, B. (eds.) ISC 2006. LNCS, vol. 4176, pp. 233–245. Springer, Heidelberg (2006)
- [Bra00] Brands, S.: Rethinking public key infrastructures and digital certificates: Building in privacy, The MIT Press,(2000) Available for download at [http://www.credentica.com/the\\_mit\\_pressbook.html](http://www.credentica.com/the_mit_pressbook.html)
- [Cha81] Chaum, D.: Untraceable electronic mail, return addresses, and digital pseudonyms. *Communications of the ACM* 24(2), 84–90 (1981)
- [Cha85] Chaum, D.: Security without identification: Transaction systems to make big brother obsolete. *Communications of the ACM* 28(10), 1030–1044 (1985)
- [CL02] Camenisch, J., Lysyanskaya, A.: Efficient non-transferable anonymous multi-show credential system with optional anonymity revocation. In: Pfitzmann, B. (ed.) EUROCRYPT 2001. LNCS, vol. 2045, pp. 93–118. Springer, Heidelberg (2001)
- [CL04] Camenisch, J., Lysyanskaya, A.: Signature schemes and anonymous credentials from bilinear maps. In: Franklin, M. (ed.) CRYPTO 2004. LNCS, vol. 3152, pp. 56–72. Springer, Heidelberg (2004)
- [CMO00] Di Crescenzo, G., Malkin, T., Ostrovsky, R.: Single database private information retrieval implies oblivious transfer. In: Preneel, B. (ed.) EUROCRYPT 2000. LNCS, vol. 1807, pp. 122–138. Springer, Heidelberg (2000)
- [CP92] Chaum, D., Pedersen, T.P.: Wallet databases with observers. In: Brickell, E.F. (ed.) CRYPTO 1992. LNCS, vol. 740, pp. 89–105. Springer, Heidelberg (1993)
- [GIKM98] Gertner, Y., Ishai, Y., Kushilevitz, E., Malkin, T.: Protecting data privacy in private information retrieval schemes. In: STOC '98: Proceedings of the thirtieth annual ACM symposium on Theory of computing, pp. 151–160. ACM Press, New York (1998)
- [GMM06] Golle, P., McSherry, F., Mironov, I.: Data collection with self-enforcing privacy. In: ACM Conference on Computer and Communications Security—ACM CCS 2006, pp. 69–78. ACM, New York (2006)
- [KO97] Kushilevitz, E., Ostrovsky, R.: Replication is not needed: single database, computationally-private information retrieval. In: FOCS '97. Proceedings of the 38th Annual Symposium on Foundations of Computer Science, pp. 364–373. IEEE Computer Society, Los Alamitos (1997)
- [Lay07] Mohamed Layouni, Controlling access to personal data through Accredited Symmetrically Private Information Retrieval, Tech. Report ??, McGill University, School of Computer Science,(July 2007)
- [Lip05] Lipmaa, H.: An oblivious transfer protocol with log-squared communication. In: Zhou, J., Lopez, J., Deng, R.H., Bao, F. (eds.) ISC 2005. LNCS, vol. 3650, pp. 314–328. Springer, Heidelberg (2005)

- [OU98] Okamoto, T., Uchiyama, S.: A new public-key cryptosystem as secure as factoring. In: Nyberg, K. (ed.) EUROCRYPT 1998. LNCS, vol. 1403, pp. 308–318. Springer, Heidelberg (1998)
- [SWP00] Song, D.X., Wagner, D., Perrig, A.: Practical techniques for searches on encrypted data. In: Proceedings of the 2000 IEEE Symposium on Security and Privacy, pp. 44–55. IEEE Computer Society, Los Alamitos (2000)