# Cryptographically Enforced Conditional Access for XML

Gerome Miklau       Dan Suciu

University of Washington

{gerome, suciu}@cs.washington.edu

## Abstract

Access control for databases is typically enforced by a trusted server responsible for permitting or denying users access to the database. This server-based protection model is increasingly becoming inconvenient for web based applications. We propose encryption techniques that allow XML documents to be distributed over the web to clients for local processing while maintaining certain access controls. In particular, we focus on conditional access controls, where a user is granted access to certain data elements conditioned on the user's existing knowledge of another part of the data. We believe such access controls are important in practice, and that enforcing them cryptographically on remote instances allows for more flexible data dissemination and processing.

## 1 Introduction

An access control model is used to permit or refuse access by subjects to data objects. Subjects are users, or groups of users usually defined by name, network identification or other static affiliation. For XML, objects are documents or parts of documents defined by XPath expressions. Access control in relational database systems, and most proposed XML systems, is enforced by a server that handles all data requests and strictly controls which users can access what data. While this model is sometimes also used in Web applications, it is often too restrictive. As the following examples show, there are a number of advantages to delivering remote copies of the data to clients if access control can be maintained:

**Local data processing** A credit card company's accounts database is secret, but if a vendor presents a correct account number and a correct expiration date the company will provide the available credit amount. Vendors could benefit if they downloaded the entire database locally: transactions could be authenticated faster and without network access, or vendors could integrate the data with their own data and run complex queries, e.g. for marketing purposes. No information would be leaked if vendors only accessed accounts for which they know the account number and expiration date.

**Privacy** An information provider that sells valuable data offers different data access rights for different prices. Customer queries are answered by the provider only if the customer has sufficient access rights. But a client's queries may reveal to the provider privileged information they prefer not to disclose[1]. Remote enforcement of access control will allow the provider to deliver an entire copy of the data so that clients can process queries locally, with privacy guaranteed.

**Offline browsing** A Web vendor does not make its products database available to the public, but, of course, allows users to browse and query for specific products. The vendor would benefit from making the database available publicly since users could download the data to their laptop and shop offline, e.g. during a plane trip.

**Peer Data Mananagement Systems** In a peer-to-peer distributed database, peers contribute data as well storage and processing resources, and members of the network can execute queries over all contributed data [12, 17]. Proposed systems require that replicas of the data be placed outside the producers' secure domains, yet data producers need to retain some access control over their data while. Remote enforcement of access controls makes this possible, and encourages peers to share data.

We propose an approach for publishing XML data on the Web while controlling how the data is accessed. In particular, we propose a novel and flexible language of *conditional access rules* used to define security policy. We explain how to encrypt XML data to enforce these access controls without a trusted server, and we discuss query processing over encrypted data.

Our notion of conditional access generalizes the static categorization of subjects into authorization

---

[1]Private information retrieval, or the problem of allowing a subject to query a remote database without revealing information through his/her queries, was first addressed in [5].

classes. Subjects are not identified by user name or network identifier but by their knowledge. As a special case, access may be conditioned on knowledge of a private key or password in a conventional way. But more generally, subjects qualify for access to an object by virtue of their knowledge of the data. Conditional access rules specify what data values need to be presented by the subject before granting access to other data values. Subject authorization is therefore flexible and dynamic in a way not possible with conventional access classes. The flexibility of our conditional access rules distinguishes our work from other attempts to encrypt data and manage decryption keys [2, 15].

Our cryptographic enforcement of conditional policies is based on known techniques for encrypting relation tables[10]. However, we view it as critical that remote access control be implemented for XML documents, rather than relations, since relations are rarely exchanged as such.

Trusting the encryption mechanism to enforce access controls on remote replicas is a substantial departure from server-based enforcement. Some applications surely require a higher level of trust in the security mechanism and will need to rely on the server model. On the other hand, there are many applications where releasing proprietary data to customers and partners is very beneficial. Our approach targets these applications, and creates opportunities for free data dissemination and new processing models.

The paper is organized as follows. In section 2 we review encryption primitives and present a simple table encryption scheme which is the basis of our document encryption method. Section 3 presents *conditional access rules*. In Section 4 we show how to generate an encrypted XML instance enforcing a set of conditional access rules, and then briefly discuss querying such an instance in Section 5. Overall security is reviewed in Section 6. We address related work and conclude in Sections 7 and 8.

## 2 Background

We present here known techniques for encrypting a relational table to enforce certain access controls. Later we adapt table encryption to enforce conditional access to XML data. We begin with definitions of encryption primitives used below [14].

**Encryption Primitives** Let $\mathcal{M}$ be the message space of plain text strings and $\mathcal{C}$ the cipher text space. A *symmetric encryption scheme* consists (for each key $k$) of an encryption function $E_k : \mathcal{M} \rightarrow \mathcal{C}$ and a corresponding decryption function $D_k : \mathcal{C} \rightarrow \mathcal{M}$ with the property that for all $m \in \mathcal{M}$, $D_k(E_k(m)) = m$. An encryption scheme is secure if it is computationally infeasible to deduce $m$ from $E_k(m)$ without knowledge of $k$. Generally, functions $D$ and $E$ are publicly known, and the security of encryption rests in the key. We also use a *collision free one-way* function $f : \mathcal{M} \rightarrow \mathcal{C}$ such that given $m \in \mathcal{M}$, it is easy to compute $f(m)$, but given $c \in \mathcal{C}$, it is computationally infeasible to find an $m$ such that $c = f(m)$. Several symmetric encryption algorithms and candidate one-way functions are mentioned in [16]. AES is a good choice for symmetric primitives $E$ and $D$. Functions that behave like $f$ can be constructed from these, or public-key encryption techniques can be used.

**Restricting Access to Relational Tables** Suppose Alice has a binary relation $T[A, B]$. She wants to publish $T$, but also wants to restrict access so that a user needs to present an $A$ value before being allowed to retrieve corresponding $B$ values. In other words, she wants to allow $T$ to be used only in the context $\sigma_{A=a}(T)$, for some constant $a$. We denote this access control rule, $r$, as $T : (A \rightarrow B)$. The solution to this problem is described in [10], expanding on a method first proposed by Needham (see [8]). Using the encryption primitives described above, the *access controlled* table (with respect to rule $r$) is $T_r^{ac}$, defined as:

$$T_r^{ac} = \{( f(a), E_a(b) ) \mid (a, b) \in T\}$$

Alice publishes $T_r^{ac}$ (along with $f$ and $D$) instead of $T$. The tuples of $T_r^{ac}$ consist only of cipher text. Any subject can use $T_r^{ac}$ in order to access $T$ in the intended manner. For example, consider the following datalog query over table $T$:

$$Q_1(x) \quad :- \quad T(\text{"abc"}, x)$$

This query accesses the data in the way it was intended, and it can be rewritten for $T_r^{ac}$ as the following rule in datalog extended with encryption primitives:

$$Q_1'(x) \quad :- \quad T_r^{ac}(f(\text{"abc"}), v), x = D_{\text{"abc"}}(v)$$

The re-written query computes $f(\text{"abc"})$, retrieves all corresponding $B$ values $v$ from $T_r^{ac}$, and then decrypts them using "abc" as the key. For a more complex example, consider $Q_2$ below and its re-writing $Q_2'$:

$$
\begin{aligned}
Q_2(x) \quad &:- \quad T(\text{"abc"}, y), T(y, z), T(z, x) \\
Q_2'(x) \quad &:- \quad T_r^{ac}(f(\text{"abc"}), v), y = D_{\text{"abc"}}(v), \\
&\qquad T_r^{ac}(f(y), u), z = D_y(u), \\
&\qquad T_r^{ac}(f(z), w), x = D_z(w)
\end{aligned}
$$

$Q_2$ still accesses the data in the way it was intended: the user first evaluates $T(\text{"abc"}, y)$, to get a set of values for $y$, then uses each of them as a constant to

evaluate $T(y, z)$, then uses each resulting $z$ to evaluate $T(z, x)$. The rewritten query $Q_2'$ shows that it is possible to express $Q_2$ on $T_r^{ac}$, by a judicious sequence of applications of $f$ and $D$.

On the other hand it is difficult for an adversary, Mallory, to use $T_r^{ac}$ in order to execute on $T$ queries that use different access patterns. For example the following queries are hard to compute without knowing values for $x$:

$$\begin{aligned} Q_3(x, y) &:- \quad T(x, y) \\ Q_4(x) &:- \quad T(x, \text{``abc''}) \end{aligned}$$

There are two known attacks to this scheme: the dictionary attack, and the guessing attack. The first applies when Mallory can generate the entire message space of $f$ which is determined by the domain of attribute $A$. For instance, when $A$ represents credit card accounts, then Mallory can enumerate all account numbers on a powerful computer and invert $f$ by brute force. The usual protection against this attack is to increase the size of the message space: take $A$ to be the account number *and* the expiration date. The second attack is the guessing attack: Mallory can simply probe random account numbers and dates, and discover $x$-values with some probability. We discuss the security of encryption further in Section 6.

# 3  Conditional Access Rules

Next we describe a flexible language of rules defining conditional access to an XML document, hence *conditional access rules* (CARs). For a path expression $P$, and an XML tree instance $D$, we denote by $eval(P, D)$ the node set resulting from evaluating $P$ starting from the root. A CAR $r$ has the following form:

$$r ::= C : (\{\bar{B}\} \to \{\bar{F}\})$$

Here $C$ is a single XPath expression specifying a *context*, $\bar{B}$ and $\bar{F}$ are sets of XPath expression defining, respectively, a set of required *bindings*, and a set of retrieved, or *free values*. $C$ is an absolute XPath expression (i.e. it starts from the document root); the expressions in $\bar{B}$ and $\bar{F}$ are relative, starting from a context node $c \in eval(C, D)$. Users specify several rules to control access to an XML document.

The intuition behind a rule is that the subtree rooted at a context node $c \in eval(C, D)$ will be hidden from the user: access is allowed only by specifying values for $\bar{B}$ (a tuple of values), in which case corresponding subtree(s) $\bar{F}$ will be accessible. All XPath expressions in $\bar{B}$ are required to return atomic values, i.e. they have to end in `text()` or `data()` or in an attribute, and we assume the set $\{\bar{B}\}$ to be ordered. Figure 1 illustrates the effect of a CAR.
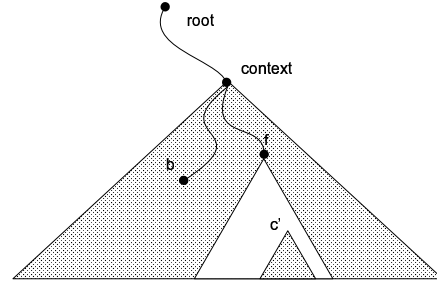


Figure 1: Illustration of the effect of a CAR. Within context node $c$, if bound value $b$ is provided, then the subtree rooted at free value $f$ is accessible, except for an inner context $c'$ that may be specified by another rule.

**Semantics**  Let $\bar{R}$ be a set of conditional access rules, and let $D$ be an XML document tree. We denote with $x \preceq y$ the fact that node $y$ is a descendant-or-self of node $x$. In a related manner we define the "allowed" descendants to be those descendants that are not beneath any other contained context node of another rule:

$$\begin{aligned} \text{allowed-desc}(z) = \{x \mid z \preceq x, \neg(\exists r \in \bar{R}, \\ r = c'(\bar{B}' \to \bar{F}'), y \in eval(c', D), z \prec y \preceq x)\} \end{aligned}$$

The semantics for CARs defines a function $access_{\bar{R}} : \mathcal{P}(Dom) \to \mathcal{P}(Dom \cup nodes(D))$, where $\mathcal{P}(S)$ denotes the set of finite subsets of $S$, with the following meaning: if the user provides values $v_1, \ldots, v_k$, then the user is allowed to see node or value $x$ if $x \in access_{\bar{R}}(\{v_1, \ldots, v_k\})$. Given $x \in nodes(D)$ we denote the value of $x$ with $val(x)$, and given $V \in Dom \cup nodes(D)$ we denote $Val(V) = \{v \mid v \in V \cap Dom\} \cup \{val(x) \mid x \in V \cap nodes(D)\}$. Then $access_{\bar{R}}$ is defined as follows:

$$\begin{aligned} access_{\bar{R}}^{(0)}(V) &= \quad V \cup \text{allowed-desc}(D_{root}) \\ access_{\bar{R}}^{(n+1)}(V) &= \quad access_{\bar{R}}^{(n)}(V) \quad \cup \\ &\quad \{x \mid x \in \text{allowed-desc}(z), \\ &\quad z \in eval(c[B_1 = v_1, \ldots, B_k = v_k]/f_j, D), \\ &\quad c : (\{B_1, \ldots, B_k\} \to \{\ldots F_j \ldots\}) \in \bar{R}, \\ &\quad v_1, \ldots, v_k \in Val(access_{\bar{R}}^{(n)}(V)) \ \} \\ access_{\bar{R}}(V) &= \quad \bigcup_{n \geq 0} access_{\bar{R}}^{(n)}(V) \end{aligned}$$

The function $access_{\bar{R}}^{(n)}(V)$ computes the nodes that the user can obtain after $n$ probes to the encrypted data: to compute some node $x$ in $access_{\bar{R}}^{(n+1)}(V)$ the

user starts by picking some values $v_1, \ldots, v_k$ already in $access_{\bar{R}}^{(n)}$ (recall that we blur the distinction between leaf nodes and their values), uses them to bind the expressions $B_1, \ldots, B_k$ in the context, then picks some node $z$ in the result of some $F_j$: the node $x$ is any descendant of $z$, except if it is hidden from the user because some other context $c'$ disallows access to it. Notice that $access_{\bar{R}}^{(0)}$ already includes all nodes in the document that are not explicitly hidden by a context.

We now provide a number of examples, based on the data pictured in Figure 2, to illustrate the flexibility of conditional access rules:

**Example 3.1** Consider the set of CARs $\bar{R}$ containing rules:

$/hospital/patientrecords/patient :$
  $(\{patientid\} \to \{pers/name, pers/address\})$
$/hospital/patientrecords/patient :$
  $(\{personal/name\} \to \{med/room, med/floor\})$

The first rule says that within each patient context, name and address are accessible conditioned on knowledge of the corresponding patient-id. The second rule says that within each patient context, the room and floor are accessible conditioned on knowledge of the corresponding patient name. If patient-id is a key for name, and name is a key for the medical fields, then together these rules imply that a patientid is sufficient to access name, address, room, and floor. All data not within a patient context is accessible.

**Example 3.2** Suppose the set of CARs $\bar{R}$ contains just the rule $/ : (\{\} \to \{.\})$. Here $C = /$, so $D_{root}$ is the only context node for this set of CARs. $B$ is the empty-set, which requires no knowledge for access to $F = \texttt{self}$ evaluated from the context, which is $D_{root}$. $access_{\bar{R}}(\{\})$ contains all descendants of $D_{root}$ since there are no other contexts to consider. Thus this rule makes the entire document accessible with no conditions.

**Example 3.3** Consider the set of CARs $\bar{R}$ containing rules:

$$/ \;\; : \;\; (\{\} \to \{.\})$$
$/hospital/patientrecords/patient \;\; :$
  $(\{patientid, medical/diagnosis\} \to \{.\})$
$$//password \;\; : \;\; (\{\} \to \{\})$$

These rules make the entire XML file accessible except (1) access to patient data is conditioned on knowing a patient's id and diagnosis, and (2) no access is given to any passwords. The last rule makes every password

element a context, within which the set of free values is empty, thus protecting all passwords, regardless of other CARs in $\bar{R}$.

**Example 3.4** Consider the set of CARs $\bar{R}$ containing the single rule:

$/hospital \;\; : \;\; ( staff/physicians/physician/$
    $\{name, password, physid = \$PID\} \to$
    $\{patrecords/patient[med/physid = \$PID]\} )$

This example shows how conventional password-authenticated access can be expressed using CARs when a password, private key, or other identifying information is included in the data. When such private information occurs in the data as bound values, it will be protected since it is within the encrypted context. The rule says that given a physician name, password, and id, the entire patient element is accessible for patients under the physician's care. (The $\$PID$ notation is a shorthand requiring matching physician id's in the bound and free parts of the rule.)

# 4   Encrypting XML

Given an XML document $D$ and a set $\bar{R}$ of conditional access rules, we generate an encrypted document $D_{\bar{R}}^{ac}$ that enforces access according to $\bar{R}$. Our strategy is to apply table encryption, as described in Section 2, to binary tables that are implied by each CAR within a context. If $r = C : (\bar{B} \to \bar{F})$ is a CAR in $\bar{R}$, and $x$ is a context node in $eval(C, D)$, then let $T_r(x)$ be the binary table which pairs the concatenation of bound data values with the concatenation of free data values. That is, $T_r(x)$ is defined by:

$$\{(b, f) \mid b = b_1.b_2 \ldots b_k, \; f = f_1.f_2 \ldots f_l,$$
$$b_i \in eval(C/B_i, D), i \in [1..k],$$
$$f_i \in eval(C/F_j, D), j \in [1..l] \quad \}$$

To construct $D_{\bar{R}}^{ac}$ we begin with $D$. Nodes that are not a context for any rule appear unchanged in $D_{\bar{R}}^{ac}$. But for other nodes we consider each rule $r$ for which $x$ is a context, compute $T_r(x)$ and replace $x$ and its subtree with the collection of tables $T_r(x)$ for each $r \in \bar{R}$. The tables are represented as XML in a standard way as a list of row elements each containing two column elements (since the $T_r(x)$ are binary). These replaced elements also need to contain some metadata describing the paths to the encoded bound and free values. A similar construction can be recursively applied to nested contexts; we omit the details from this abstract.

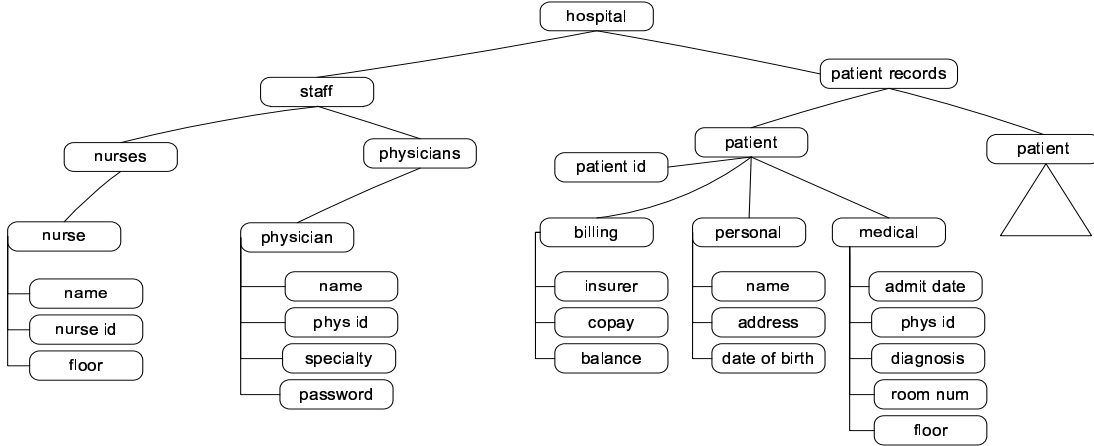Admittedly, this is a rather naive encryption scheme. The resulting $D_{\bar{R}}^{ac}$ may be large compared

4

Figure 2: Tree view of XML data describing hospital staff and patients (data values omitted).

with $D$, and grows with the number of conditional access rules in $\bar{R}$. Furthermore, there is an opportunity for optimizing the set of CARs to produce smaller encrypted instances. Given a set of CARs $\bar{R}$, we would like to find an equivalent[2] set $\bar{R}'$ that is minimal.

Surprisingly, we can use the theory of functional dependencies to aid minimization of a set of access rules, at least in the case of relational data. For example, given a ternary table $T[A, B, C]$ with the set of relational access rules (as in Section 2) $T : (A \to B), T : (B \to C), T : (A \to C)$ the naive encrypted version will represent consist of three binary relations, encoding $\Pi_{AB}(R), \Pi_{BC}(R), \Pi_{AC}(R)$. The first two relations, however, are sufficient. This happens because the given set of access rules is equivalent to the following: $T : (A \to B), T : (B \to C)$. For a set of relational access rules $\bar{S}$, denote $\bar{S}^{fd}$ the functional dependencies obtained from interpreting $\to$ as being a functional dependency rather than access control. For a relational schema, $T : A \to B$ has the standard meaning: $A$ functionally determines $B$ in $T$.

**Theorem 4.1** *Two sets of relational access rules $\bar{S}_1, \bar{S}_2$ are equivalent iff $\bar{S}_1^{fd}$ is equivalent to $\bar{S}_2^{fd}$.*

For XML, we can interpret a CAR $C : (B \to F)$ as a functional dependency according to the definition given in [3]. However, the above theorem does not hold for XML documents. We leave the interesting connection between functional dependencies and conditional access rules, along with the general CAR minimization problem, as a compelling direction for future work.

---

[2] Two CAR sets $\bar{R}_1, \bar{R}_2$ are equivalent if for any XML document $D$, $access_{R_1}(D) = access_{R_2}(D)$.

# 5 Query Processing of Encrypted XML

Since an encrypted, access-controlled document $D^{ac}$ may contain large portions of the original document unencrypted, it makes sense to use existing tools and technologies to query $D^{ac}$. We consider XQuery[4] a suitable general-purpose query language. For a query $q$ and an XML document $D$ we denote by $q(D)$ the answer of $q$ evaluated on $D$. The query rewriting problem is to find a new query $q'$ such that $q(D) = q'(D^{ac})$. The new query $q'$ will consist of XQuery syntax augmented with operators for the one-way function $f$ and the decryption function $D$ needed to process the data. In particular, the parts of the query that attempt to access protected contexts need to be replaced with routines that extract data from the encrypted tables represented as XML beneath the context. These routines are closely related to the rewriting examples provided in Section 2. A full treatment of the re-writing problem is omitted from this abstract.

# 6 Security

As we mentioned in Section 2, there are two explicit attacks on table encryption (dictionary and guessing), and these vulnerabilities are inherited by our encrypted instances. Both attacks depend on the size of the message space, determined by the product of the domain sizes of the bound attributes in a CAR. The main protection against a guessing attack is to decrease the probability of success by increasing the domain. It is worth noting that a secure server allowing similar access is also vulnerable to a guessing attack, however this can be dealt with by monitor-

ing queries and delaying the response between unsuccessful probes. In a similar manner, the computation speed of the one-way function will slow down probes to the encrypted instance but will also impede query processing.

In addition, the access controlled relational table $T^{ac}$ leaks information about $T$ in a subtle way. From $T^{ac}$ the subject can compute the number of tuples in $T$, and the number of distinct values in column $A$. The latter form of information leakage can be avoided by using an improved table encryption scheme omitted for ease of exposition.

# 7 Related Work

There are a number of recently proposed access control models for XML. In [6, 7, 11] XPath expressions identify fine-grained authorization objects and a server negotiates access and delivers a pruned document. A similar framework is proposed in [2] and then extended to offer remote enforcement by encrypting elements and transferring encrypted keys to users. Our work generalizes this by making access conditioned on encrypted keys, or on other knowledge of the data. The authors of [13] describe a model where provisional access policies can be defined. Such a policy grants a subject access to data provided a further action or qualification is satisfied (e.g., a subject's activity may be logged, or a password is provided). These provisions are related to our notion of conditional access control although they work within the framework of a trusted security processor. A framework for element-level encryption of XML documents has been proposed by the W3C [9]. The IBM Security Suite[1] implements the W3C XML encryption recommendation.

# 8 Conclusions and Future Work

We have proposed techniques for enforcing access control without relying on a secure server, with many practical applications. Removing the server from access control enables offline query data processing, anonymous evaluation of sensitive queries, and new data sharing possibilities. In addition, remotely controlled access can still be useful in a setting where data is served to clients because server contractors may themselves be untrusted. Further, we hope our conditional access rules would make it easy for individuals to publish sensitive data in the absence of sophisticated security systems. We have identified the following key problems that require further attention: (1) minimizing a set of conditional access rules to improve the efficiency of encryption and aid users in policy definition; (2) re-writing queries over encrypted instances; and (3) quantifying the security of encrypted

instances in terms of the attribute domains and properties of the encryption primitives.

# References

[1] AlphaWorks. Xml security suite. www.alphaworks.ibm.com/tech/xmlsecuritysuite, November 2001.

[2] E. Bertino, S. Castano, and E. Ferrari. Securing xml documents with author-x. *IEEE Internet Computing*, May/June 2001.

[3] P. Buneman, S. Davidson, W. Fan, C. Hara, and W. Tan. Reasoning about keys for xml. *Workshop of Databases and Programming Languages (DBLP 2001)*, 2001.

[4] D. Chamberlin, J. Clark, D. Florescu, J. Robie, J. Simeon, and M. Stefanascu. XQuery 1.0: An XML query language. http://www.w3.org/TR/xquery/, 07 June 2001. W3C working draft.

[5] B. Chor, O. Goldreich, E. Kushilevitz, and M. Sudan. Private information retrieval. In *Foundations of Computer Science (FOCS)*, 1995.

[6] E. Damiani, S. D. C. di Vimercati, S. Paraboschi, and P. Samarati. Controlling access to xml documents. *IEEE Internet Computing*, 5(6):18–28, November/December 2001.

[7] E. Damiani, S. D. C. di Vimercati, S. Paraboschi, and P. Samarati. A fine-grained access control system for xml documents. *ACM Transactions on Information and System Security (TISSEC)*, 2002. To appear.

[8] D. Denning. *Cryptography and Data Security*. Addison-Wesley Publishing Co., 1982.

[9] D. Eastlake and J. Reagle. Xml encryption syntax and processing. http://www.w3.org/TR/xmlenc-core, 4 March 2002. World Wide Web Consortium (W3C) Working Draft.

[10] J. Feigenbaum, M. Y. Liberman, and R. N. Wright. Cryptographic protection of databases and software. In *Distributed Computing and Cryptography*, pages 161–172, 1991.

[11] A. Gabillon and E. Bruno. Regulating access to xml documents. *Proc. of the 15th Annual IFIP WG 11.3 Working Conference on Database and Application Security*, July 2001.

[12] S. Gribble, A. Halevy, Z. Ives, M. Rodrig, and D. Suciu. What can peer-to-peer do for databases, and vice versa? WebDB Workshop on Databases and the Web, June 2001.

[13] M. Kudo and S. Hada. Xml document security based on provisional authorization. *Computer and Communication Security (CCS)*, November 2000.

[14] A. J. Menezes, P. C. van Oorschot, and S. A. Vanstone. *Handbook of Applied Cryptography*. CRC Press, 1997.

[15] J. Reagle. Xml encryption requirements. http://www.w3.org/TR/xml-encryption-req, 4 March 2002. World Wide Web Consortium (W3C) Note.

[16] B. Schneier. *Applied Cryptography, Second Edition*. John Wiley and Sons, Inc., 1996.

[17] M. Stonebraker, P. M. Aoki, W. Litwin, A. Pfeffer, A. Sah, J. Sidell, C. Staelin, and A. Yu. Mariposa: A wide-area distributed database system. *VLDB Journal: Very Large Data Bases*, 5(1):48–63, 1996.