

# XML データベースへの推論攻撃による機密情報特定可能性の形式化とある前提条件のもとでの特定可能性検証法の提案

A Formalization of the Identifiability of Secret Information by Inference Attacks on XML Databases and a Proposal of an Identifiability Verification Method under Some Conditions

高須賀 史和<sup>\*</sup> 橋本 健二<sup>\*</sup>  
石原 靖哲<sup>\*</sup> 藤原 融<sup>\*</sup>

Fumikazu TAKASUKA Kenji HASHIMOTO  
Yasunori ISHIHARA Toru FUJIWARA

推論攻撃とは、ユーザが許可されている問合せとその実行結果から、許可されていない問合せの実行結果（機密情報）を推論しようとするものである。本論文では、ユーザが許可されている複数の問合せとその結果から、機密情報の候補を絞り込む攻撃と、その攻撃により機密情報が特定される可能性を考える。そして、XML データベースへの推論攻撃による機密情報の特定可能性を形式的に検証する手法を提案する。

Inference attacks mean that an attacker tries to infer the execution result of a query unauthorized to the attacker (i.e., secret information) from the execution results of queries authorized to the attacker. In this paper, we consider attacks such that the attacker narrows the candidates for the secret information using more than one authorized queries and their execution results, and the possibility that the secret information can be identified by the attacks. Then, a method of verifying the identifiability by the inference attacks on XML databases is proposed.

## 1. まえがき

企業や組織は少なからず機密性の高い情報を保持している。そのような機密情報が格納されているデータベースシステムにおいて、セキュリティ面での管理が必ずしも正確に行われていないことが問題となっている。たとえば、複数のセ

<sup>\*</sup>学生会員 大阪大学大学院情報科学研究科 博士前期課程  
fmkz-tak@ist.osaka-u.ac.jp

<sup>\*</sup>大阪大学大学院情報科学研究科 博士後期課程  
k-hasimt@ist.osaka-u.ac.jp

<sup>\*</sup>大阪大学大学院情報科学研究科  
fishihara.fujiwara@ist.osaka-u.ac.jp

キュリティレベルのユーザをもつデータベースシステムにおいて、アクセス権の付与が適切に行われておらず、本来権限のないユーザによる機密情報へのアクセスを許可してしまっている場合がある。また、機密情報への直接のアクセスを許可していない場合でも、アクセスを許可されている情報やドメイン知識などを基に、本来権限の無いユーザがその機密情報を推論できてしまう場合がある[1]。そのような推論を行うことを推論攻撃と呼ぶ。推論攻撃によりどのような情報が得られるかは一般に自明ではないため、データベース管理者としては、データベースの運用を開始する前、あるいはユーザからの問合せに対して実行結果を返す前に、機密情報に対して推論攻撃が成功する可能性をあらかじめ把握し、その推論を制御できることが重要である。

[例1] XMLデータベースにおける推論攻撃の例を示す。Dは学生名とその成績の対応を表すXML文書であり、以下のスキーマにしたがっているとする。

科目	学生 <sup>*</sup>
学生	氏名, 点数
氏名	PCDATA
点数	PCDATA

すなわち、科目要素の子には学生要素が0個以上並び、各学生要素の子には氏名要素と点数要素が1つずつ並び、氏名要素と点数要素は文字列(PCDATA)を値としてとる。

$T_1$ を、もとの文書における出現順序を保存したまま、科目の子の学生の子の氏名だけを取り出す問合せとする。また、 $T_2$ を、もとの文書における出現順序を保存したまま、科目の子の学生の子の点数だけを取り出す問合せとする。

今、 $T_1$ および $T_2$ の実行が許可されており、実行結果 $T_1(D)$ 、 $T_2(D)$ がそれぞれ図1、2に示すとおりであったとする。 $T_1$ の定義および $T_1(D)$ の値より、Dは科目の子の学生の子の氏名として“高須賀”という値の要素と“橋本”という値の要素をその順にもつということがわかる。同様に、 $T_2$ の定義および $T_2(D)$ の値より、Dは科目の子の学生の子の点数として“90”という値の要素と“80”という値の要素をその順にもつということがわかる。そして、Dが上のスキーマにしたがっているということより、Dは図3に示す文書しかありえないと結論できる。



図1 実行結果 $T_1(D)$   
Fig.1 Execution Result  
 $T_1(D)$

図2 実行結果 $T_2(D)$   
Fig.2 Execution Result  
 $T_2(D)$

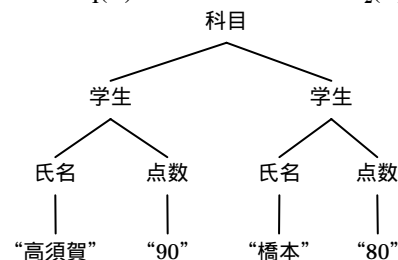


図3 Dとして推論される文書  
Fig.3 Inferred Document for D

次に、本論文で扱う、機密情報の候補を絞り込む推論攻撃の概略を示す(図4参照)。ユーザはXML文書 $D$ のある部分 $T_S(D)$ (ここで $T_S$ は問合せ)の値を知りたい。しかし、ユーザは $D$ 全体を直接見ることができない。ユーザが知ることが出来るのは、実行を許可されている問合せ $T_1, \dots, T_p$ の定義と、これらの問合せの $D$ における実行結果 $T_1(D), \dots, T_p(D)$ 、および $D$ がしたがうスキーマ $G$ である。このとき、各 $i$ ( $1 \leq i \leq p$ )について、ユーザは $T_i$ と $T_i(D)$ とから、問合せ $T_i$ の実行結果として $T_i(D)$ が得られるような $D$ を推論することにより、 $D$ の候補を絞り込むことができる。すべての $i$ について $D$ の候補となっている文書のうち、 $G$ にしたがう文書からなる集合が、 $D$ の候補集合 $I_D$ である。そして、どの候補 $D' \in I_D$ に対しても $T_S(D')$ が同じ値になるのであれば、ユーザは $T_S(D)$ の値を特定できる。

本論文では、まず、機密情報の特定可能性を形式的に定義する。問合せのクラスとしては、文献[2]で提案されているトップダウン木変換器を用いる。この木変換器はXSLTの部分クラスに相当する。次に、ある前提の下で機密情報の特定可能性を検証する手法を与える。本検証法の時間計算量は、実行が許可された問合せの個数 $p$ を定数とみなすと、決定性多項式時間である。

関連研究として、文献[3]では、各ユーザに security view を提供するアプローチを提案している。security viewは、ユーザがアクセス権をもつ情報をちょうど含んだ文書と、その文書のスキーマとから成る。問合せ言語としてはXPathを前提としている。元のスキーマからのビュースキーマの生成に関しては、推論攻撃の手がかりとなりうる要素名を付け替えたり削除したりするアルゴリズムを提案している。しかし、安全性の定義が不明確であるという問題点がある。一方、文献[4]は、推論攻撃を受けたとしても情報の漏洩なく公開できる部分文書を算出するアルゴリズムを提案している。攻撃者が推論に用いることが出来る情報は、スキーマ情報と関数従属性に関する情報である。本論文では関数従属性に関する情報を扱ってはいないが、[4]では考慮していない、複数の問合せ結果から機密情報の候補を絞り込む攻撃を扱うことができる。

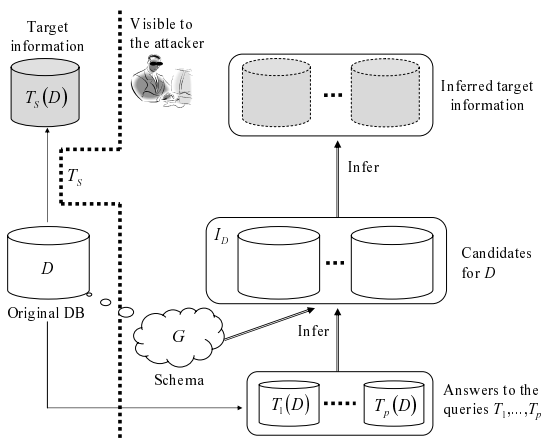


図4 機密情報の候補を絞り込む推論攻撃

Fig.4 The Inference Attack Trying to Narrow the Candidates for Secret Information

## 2. 準備

木、生垣、非決定性トップダウン木オートマトン、決定性トップダウン木変換器に関する諸定義を行う。

### 2.1 木、生垣

あるアルファベット $\Sigma$ について、 $\Sigma$ 上の(ランクなし)木の集合 $T_\Sigma$ と生垣の集合 $H_\Sigma$ を、以下を満たす最小集合と定義する:

- $H_\Sigma = T_\Sigma^*$ ,
- $a \in \Sigma, h \in H_\Sigma$ であるならば、 $a(h) \in T_\Sigma$ .

ここで $T_\Sigma^*$ は $T_\Sigma$ のKleene閉包を表す。以下、空生垣(木の空系列)を $\varepsilon$ で表す。木において、空生垣を子としても頂点をその木の葉頂点と呼ぶ。木 $t$ (あるいは生垣 $h$ )の大きさ $|t|$ (あるいは $|h|$ )はその木(あるいは生垣)の頂点の数である。

### 2.2 非決定性トップダウン木オートマトン

[定義1] 非決定性トップダウン木オートマトンは次の4つ組 $A=(Q, \Sigma, q_0, R)$ である:

- $Q$ : 状態の有限集合,
- $\Sigma$ : アルファベット,
- $q_0 \in Q$ : 初期状態,
- $R$ : 遷移規則の集合.

ただし $q \in Q, a \in \Sigma$ とし、 $e$ を $Q$ 上の言語を受理する非決定性有限オートマトンとすると、遷移規則は $(q, a) \rightarrow e$ の形式である。

以下、木オートマトン $A=(Q, \Sigma, q_0, R)$ の動作について述べる。 $t=a(t_1 \dots t_n)$ とする。ここで、 $t_1, \dots, t_n \in T_\Sigma$ である。木 $t$ の根頂点に状態 $q$ が割り当てられたとき $q(t)$ と書く。 $q_1 \dots q_n \in L(e)$ であるような $((q, a) \rightarrow e) \in R$ が存在するならば、 $A$ は $q(t)$ から $a(q_1(t_1) \dots q_n(t_n))$ へ遷移可能であると定義する。ここで $L(e)$ は非決定性有限オートマトン $e$ が受理する文字列言語である。 $q_0(t)$ から始めて、以上のような遷移をトップダウンに各部分木について繰り返すことによって、最終的に $t$ へ遷移可能であるとき、 $A$ は $t$ を受理するという。 $A$ によって受理されるすべての木の集合を $L(A)$ と書く。

$A$ の大きさは $|Q| + |\Sigma| + \sum_{q \in Q, a \in \Sigma} |\text{rhs}(q, a)|$ である。ここで、 $|\text{rhs}(q, a)|$ は $(q, a) \rightarrow e$ について非決定性有限オートマトン $e$ の大きさである。非決定性有限オートマトンの大きさも、同様に、状態数と記号数、そして各遷移規則における遷移先の状態集合の大きさの総和である。

木オートマトンはXMLスキーマの形式的モデルの1つであり、XML文書の“型”を表現するのに用いられる。ある木変換器 $T$ と木オートマトン $A_{in}$ について、 $L(A'_{out}) = \{T(t) \mid t \in L(A_{in})\}$ であるような木オートマトン $A'_{out}$ を求めることを型推論という。逆に、ある木変換器 $T$ と木オートマトン $A_{out}$ について、 $L(A'_{in}) = \{t \mid T(t) \in L(A_{out})\}$ であるような木オートマトン $A'_{in}$ を求めることを逆型推論という。

### 2.3 決定性トップダウン木変換器

$Q$ 中の記号が葉頂点のみに現れるような $\Sigma$ 上の木集合を $T_\Sigma(Q)$ と書く。 $H_\Sigma(Q)$ も同様に定義する。

[定義2] 決定性トップダウン木変換器は次の4つ組 $T=(Q, \Sigma, q_0, R)$ である:

- $Q$ : 状態の有限集合,
- $\Sigma$ : アルファベット,
- $q_0 \in Q$ : 初期状態,
- $R$ : 変換規則の集合.

ただし、 $q \in Q, \{\lambda\}, h \in H_\Sigma(Q)$ とすると、変換規則は $(q, a) \rightarrow h$ の形式である。また $q=q_0$ のとき $h$ は $T_\Sigma(Q) \setminus Q$ の要素でなければならない。ここで $\lambda$ は空記号である。また、決定性トップダウン木変換器では、同じ左辺を持つ変換規則はちょうど1つである。また、 $((q, a) \rightarrow h) \in R$ であるすべての $q \in Q$ とすべての $a \in \Sigma$ について、 $((q, a) \rightarrow h) \notin R$ であり、 $((q, a) \rightarrow h) \in R$ であるすべ

ての $q \in Q$ とすべての $a \in \Sigma$ について、 $((q, a) \ h) \in R$ である。決定性トップダウン木変換器 $T=(Q, \Sigma, q_0, R)$ の動作について述べる。木 $t = a(t_1 \dots t_n)$ は、状態 $q \in Q$ で $T$ によって次のような $T^q(t)$ へ変換される:

- ある $h$ に対して $((q, \ ) \ h) \in R$ ならば、 $T^q(t)$ は、 $h$ 中のすべての状態記号 $p$ を $T^p(t)$ に置換して得られる生垣と等しい。
- ある $h$ に対して $((q, a) \ h) \in R$  (ただし、 $a \in \Sigma$ )ならば、 $T^q(t)$ は、 $h$ 中のすべての状態記号 $p$ を生垣 $T^p(t_1) \dots T^p(t_n)$ に置換して得られる生垣と等しい。特に $t=a(\varepsilon)$ のときは、 $T^q(t)$ は、 $h$ 中のすべての状態記号 $p$ を $\varepsilon$ に置換して得られる生垣と等しい。

$t$ の $T$ による変換 $T(t)$ を $T(t) = T^{q_0}(t)$ と定義する。 $T$ の大きさは $|Q| + |\Sigma| + \sum_{q \in Q, a \in \Sigma} |\text{rhs}(q, a)|$ である。ここで $\text{rhs}(q, a)$ は $((q, a) \ h) \in R$ について $|h|$ である。初期状態に関する変換規則の制約により、この変換器の出力が木であることが保証される。また変換規則の中で、右辺に状態が2つ以上存在する規則は、規則を適用している頂点以下の部分木をコピーするので copying rule と呼び、右辺にアルファベットを含まない規則は、規則を適用している頂点を除去する変換を行うので deletion rule と呼ぶ。そして、copying rule を持たない木変換器を non-copying, deletion rule を持たない木変換器を non-deleting と呼ぶ。

### 3. 機密情報の特定可能性

本節では、機密情報の特定可能性を定義する。ユーザに与えられる情報は以下の通りである。

- $T_1, \dots, T_p$ : ユーザが実行を許可されている問合せ。
- $T_1(D), \dots, T_p(D)$ :  $D$ に対する問合せ $T_1, \dots, T_p$ の実行結果。
- $A_G$ :  $D$ がしたがうスキーマ (を表す木オートマトン)。
- $T_S$ :  $D$ 中の機密情報を返す問合せ。

そしてこれらから定まる以下の集合

$I_S = \{T_S(D') \mid D' \in L(A_G), T_1(D') = T_1(D), \dots, T_p(D') = T_p(D)\}$ が要素を1つしかもたないとき、 $D$ の機密情報 $T_S(D)$ は特定可能であるという。

たとえば[例 1]では、スキーマ $A_G$ にしたがうXML文書の中で問合せ $T_1$ の実行結果として図1に示す $T_1(D)$ 、 $T_2$ の実行結果として図2に示す $T_2(D)$ が得られるような $D$ の候補 $D'$ は、 $T_1, T_2$ の定義より図3に示すものただ1つであると結論付けられる。そして、その文書における機密情報(この場合、学生の氏名と点数の組)を返す問合せ $T_S$ (この場合は $T_S(D') = D'$ )を実行した結果、得られる集合 $I_S$ の要素もただ1つであるため、 $D$ の機密情報 $T_S(D)$ も特定可能である。

### 4. 特定可能性検証法

本節では、本論文において提案する機密情報の特定可能性検証法について述べる。

機密情報の特定可能性の検証は次のように行う。まず、ユーザが実行を許可されている問合せの実行結果と $D$ がしたがうスキーマから $D$ の候補を求める。次に、得られた $D$ の候補と、 $D$ の機密情報を返す問合せから $D$ 中の機密情報の候補を求め、そして、その候補が1つであるかどうか、すなわち、機密情報が特定可能であるかを判定する。

以降4.1節では本検証を行うための前提条件について、4.2節では本検証の手順、そして4.3節では本検証の時間計算量の評価について述べる。

#### 4.1 特定可能性検証のための前提条件

ここでは検証を行うための前提条件について述べる。本論

文では、まずユーザが実行を許可されている各問合せ $T_i$ として、non-copyingかつnon-deletingな木変換器 $T_{i1}$ とnon-copyingな木変換器 $T_{i2}$ の合成により定義されている問合せを考える。 $T_{i1}$ は、コピーや削除以外の変換を行いつつ、後に $T_{i2}$ によって削除されるべき頂点のラベルを新たなラベル#に変換する機能を持つ。そして $T_{i2}$ は#とラベル付けされた頂点を削除するだけの機能を持つ木変換器である。また、 $D$ 中の機密情報を求める問合せ $T_S$ として、non-copyingな木変換器により定義されている問合せを考える。

#### 4.2 特定可能性検証の手順

ここでは特定可能性検証の手順について述べる。本検証法は以下の5つのステップにより構成される(図5参照)。

- (1) 各問合せの実行結果 $T_i(D)$ から、 $T_{i2}(D)$ の出力として $T_{i1}(D)$ が得られるような $T_{i2}$ への入力木の候補集合を受理するオートマトン $A_{i1}^d$ を求める。
- (2) ステップ(1)で得られた各木オートマトン $A_{i1}^d$ と木変換器 $T_{i1}$ から、逆型推論を用いて、元のXML文書 $D$ の候補を受理する木オートマトン $A_{iD}^d$ を構成する。
- (3) ステップ(2)までで求められたすべての $A_{iD}^d$ と $D$ がしたがっているスキーマを表す木オートマトン $A_G$ から、 $D$ の候補を受理する木オートマトン $A_D$ を求める。
- (4) ステップ(3)で求められた $A_D$ と機密情報を求める木変換器 $T_S$ から型推論を用いて、XML文書 $D$ 中の機密情報 $T_S(D)$ の候補を受理する木オートマトン $A_S$ を求める。
- (5) ステップ(4)で求められた $A_S$ が受理する木がただ1つであるかを判定する。

これにより機密情報 $T_S(D)$ が特定可能であるかどうかを判定する。

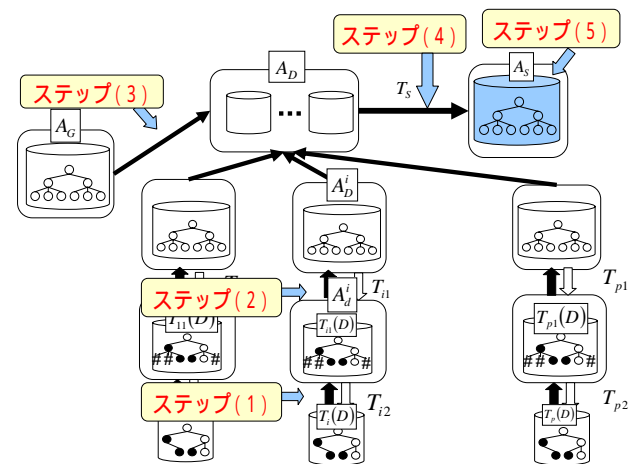


図5 機密情報の特定可能性検証  
Fig.5 Identifiability Verification of Secret Information

#### 4.3 時間計算量の評価

ここでは、本検証法の時間計算量について述べる。まずステップ(1)の計算については、この計算量は $|T_i(D)|$ についてのLOGSPACEである[6]。またここで求められる $A_{i1}^d$ の大きさは、 $|T_i(D)|$ の多項式の大きさとなる。

次に、ステップ(2)の計算は、non-copyingかつnon-deletingの木変換器 $T_{i1}$ の各変換規則につき、その右辺の記号を左辺に持つ $A_{i1}^d$ の遷移規則を選び、それらから新たに

遷移規則を作るという作業を繰り返すことから、時間計算量は  $O(|T_{il}| \times |A_{il}^i|)$  となる。

ステップ(3)の計算の時間計算量は、 $O(|A_{il}^1| \times \dots \times |A_{il}^p|)$  となる。 $p$  を定数とすると多項式時間となる。 $|A_{il}^p|$  についても、状態集合の大きさが各オートマトンの状態集合の大きさの積となるので、これも時間計算量と同様、 $p$  を定数とすると多項式の大きさとなる。

ステップ(4)の計算は、ステップ(2)を求める計算と同様、各変換規則と遷移規則から新たな遷移規則を作るという作業になるため、この時間計算量は  $O(|T_{il}| \times |A_{il}^p|)$  となる。

最後にステップ(5)の計算について考える。このステップでは、まずトップダウンの木オートマトンを等価なボトムアップの木オートマトンに書き換える。この時間計算量は木オートマトンのサイズについて線形時間である。次に、いま得られたオートマトンを非決定性(ボトムアップ、ランクなし)のものから、二分木オートマトンへの変換を行う。この実行時間は多項式時間であることがわかっている[7]。そしてこのボトムアップ二分木(ランクあり)オートマトンが受理する木の集合が要素をただ1つだけもつかどうかを判定するが、この計算も実行時間は多項式時間である[5]。よってステップ(5)の計算は  $|A_{il}^p|$  についての多項式時間である。

以上より、本検証法の時間計算量は  $A_{il}^p$  を求める際に交わりを求める木オートマトンの数(すなわちユーザが実行を許可されている問合せの個数  $p$ ) を定数とみなしたとき、多項式時間である。

## 5. あとがき

本論文ではXMLデータベースにおける推論攻撃による機密情報の特定可能性を形式的に定義し、ある前提条件の下でそれを検証する手法を提案した。本検証法の時間計算量は、ユーザが実行を許可されている問合せの個数を定数とみなしたとき、決定性多項式時間である。

今後の課題としては、まず、検証のための前提条件の緩和が挙げられる。本検証のための前提条件では変換規則への制限が多い。たとえば、変換規則の右辺に生垣や枝分かれのある木を許していないため、問合せに対して新たな追加情報を加えて返すような動作ができない。そこで、この制限を緩めることについて検討をはじめている。

また、問合せクラスの拡張も課題として挙げられる。たとえば「A という子要素をもつときかつそのときのみ、親要素を削除する」という動作ができない。そこで、子や孫への先読み機能をもつ木変換器に拡張することについて検討をはじめている。

次に、本論文で取り上げている機密情報の特定可能性とデータベースの安全性についての議論があげられる。機密情報が1つに絞り込まれる(特定される)可能性を本研究では考えているが、一方、機密情報が1つに絞り込まれない場合においても、以下のような両極端な状況が存在しうる。

- 機密情報の候補の個数が数個にしかならない(特定されないまでも攻撃の効果が非常に高い)
- 機密情報の候補の個数が非常に多い(攻撃の効果が低い、あるいは非常に低い)

実用上、これら両極端な状況は区別されるべきケースが多いと考えられ、今後、この間に安全かそうでないかの線引きを行うことが課題となる。そこで、以下のような問題について検討をはじめている。

- 機密情報の候補が有限個になるかを検証する。

- ある整数のパラメータ  $k$  を与え、機密情報の候補が  $k$  個以下になるかどうかを検証する。

他に、推論に用いることのできる情報として、[4]のような関数従属性を考慮する必要があると考えられる。また、本論文はXML文書  $D$  が具体的に与えられている場合での安全性検証を扱っているが、静的な特定可能性検証、すなわち、ユーザに実行を許可する問合せと、スキーマ情報から、機密情報が特定可能であるような文書が存在するか否かを検証する問題も、実用上極めて重要な課題である。

## [文献]

- [1] C. P. Pfleeger and S. L. Pfleeger: "Security in Computing," 3rd Ed., Prentice Hall (2002).
- [2] W. Martens and F. Neven: "Frontiers of Tractability for Typechecking Top-Down XML Transformations," Proceedings of the 23rd ACM SIGMOD SIGACT SIGART symposium on Principles of database systems, pp. 23-34 (2005).
- [3] W. Fan, C.-Y. Chan, and M. Garofalakis: "Secure XML Querying with Security Views," Proceedings of the 23rd SIGMOD International Conference on Management of Data, pp. 587-598 (2004).
- [4] X. Yang and C. Li: "Secure XML Publishing without Information Leakage in the Presence of Data Inference," Proceedings of the 30th VLDB Conference, pp. 96-107 (2004).
- [5] H. Comon, M. Dauchet, R. Gilleron, F. Jacquemard, D. Lugiez, S. Tison, and M. Tommasi: "Tree Automata Techniques and Applications," <http://www.grappa.univ-lille3.fr/tata/>.
- [6] W. Martens and F. Neven: "On the Complexity of Typechecking Top-Down XML Transformations," Theoretical Computer Science, Volume 336, Issue 1, pp. 153-180 (2005).
- [7] F. Neven: "Automata theory for XML researchers," ACM SIGMOD Record Volume 31, Issue 3, pp. 39-46 (2002).

## 高須賀 史和 Fumikazu TAKASUKA

大阪大学大学院情報科学研究科博士前期課程在学中。2005 神戸大学工学部情報知能工学科卒業。XML データベースやデータベースセキュリティに関心をもつ。

## 橋本 健二 Kenji HASHIMOTO

大阪大学大学院情報科学研究科博士後期課程在学中。2006 大阪大学大学院情報科学研究科博士前期課程修了。XML データベースやデータベースセキュリティに関心をもつ。

## 石原 靖哲 Yasunori ISHIHARA

大阪大学大学院情報科学研究科助教授。1990 大阪大学基礎工学部情報工学科卒。1992 同大学院基礎工学研究科博士前期課程了。1994 同大学院基礎工学研究科博士後期課程退学。同年より奈良先端科学技術大学院大学情報科学研究科助手。博士(工学)。XML データベースやデータベースセキュリティに関心をもつ。ACM, IEEE, 情報処理学会, 電子情報通信学会各会員。

## 藤原 融 Toru FUJIWARA

大阪大学大学院情報科学研究科教授。1981 大阪大学基礎工学部情報工学科卒。1986 同大学院博士課程了。博士(工学)。情報セキュリティ, 符号理論等の研究に従事。ACM, IEEE, 情報処理学会, 電子情報通信学会各会員。