

# On XML Integrity Constraints in the Presence of DTDs

Wenfei Fan  
Bell Laboratories  
and  
Leonid Libkin  
University of Toronto

---

The paper investigates XML document specifications with DTDs and integrity constraints, such as keys and foreign keys. We study the consistency problem of checking whether a given specification is meaningful: that is, whether there exists an XML document that both conforms to the DTD and satisfies the constraints. We show that DTDs interact with constraints in a highly intricate way and as a result, the consistency problem in general is undecidable. When it comes to unary keys and foreign keys, the consistency problem is shown to be NP-complete. This is done by coding DTDs and integrity constraints with linear constraints on the integers. We consider the variations of the problem (by both restricting and enlarging the class of constraints), and identify a number of tractable cases, as well as a number of additional NP-complete ones. By incorporating negations of constraints, we establish complexity bounds on the implication problem, which is shown to be coNP-complete for unary keys and foreign keys.

Categories and Subject Descriptors: H.2.1 [**Database Management**]: Data models; F.4.3 [**Mathematical Logic and Formal Languages**]: Decision problems; I.7.2 [**Document Preparation**]: Markup languages

General Terms: Algorithms, Design, Languages, Theory

Additional Key Words and Phrases: Integrity Constraints, DTDs, XML, Consistency, Implication

---

---

Name: Wenfei Fan

Address: 600 Mountain Avenue, Murray Hill, NJ 07974, USA. Email: wenfei@research.bell-labs.com. On leave from Temple University. Supported in part by NSF Career Award IIS-0093168.

Name: Leonid Libkin

Affiliation: Department of Computer Science

Address: Department of Computer Science, University of Toronto, Toronto, Ontario, M5S 3H5, Canada. Email: libkin@cs.toronto.edu. Supported in part by National Sciences and Engineering Research Council of Canada.

Research affiliation: Bell Laboratories.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or direct commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works, requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept, ACM Inc., 1515 Broadway, New York, NY 10036 USA, fax +1 (212) 869-0481, or [permissions@acm.org](mailto:permissions@acm.org).

## 1. INTRODUCTION

Although a number of dependency formalisms were developed for relational databases, functional and inclusion dependencies are the ones used most often. More precisely, only two subclasses of functional and inclusion dependencies, namely, keys and foreign keys, are commonly found in practice. Both are fundamental to conceptual database design, and are supported by the SQL standard [30]. They provide a mechanism by which one can uniquely identify a tuple in a relation and refer to a tuple from another relation. They have proved useful in update anomaly prevention, query optimization and index design [1; 37].

XML (eXtensible Markup Language [6]) has become the prime standard for data exchange on the Web. XML data typically originates in databases. If XML is to represent data currently residing in databases, it should support keys and foreign keys, which are an essential part of the semantics of the data. A number of key and foreign key specifications have been proposed for XML, e.g., the XML standard (DTD) [6], XML Data [27] and XML Schema [36]. Keys and foreign keys for XML are important in, among other things, query optimization [34], data integration [21], and in data transformations between XML and database formats [28].

XML data usually comes with a DTD<sup>1</sup> that specifies how a document is organized. Thus, a specification of an XML document may consist of both a DTD and a set of integrity constraints, such as keys and foreign keys. A legitimate question then is whether such a specification is *consistent*, or meaningful: that is, whether there exists a (finite) XML document that both satisfies the constraints and conforms to the DTD.

In the relational database setting, such a question would have a trivial answer: one can write arbitrary (**primary**) **key** and **foreign key** specifications in SQL, without worrying about consistency. However, DTDs (and other schema specifications for XML) are more complex than relational schema: in fact, XML documents are typically modeled as node-labeled trees, e.g., in XSL [15], XQL [35], XML Schema [36], XPath [16] and DOM [3]. Consequently, DTDs may interact with keys and foreign keys in a rather nontrivial way, as will be seen shortly. Thus, we shall study the following family of problems, where  $\mathcal{C}$  ranges over classes of integrity constraints:

### XML SPECIFICATION CONSISTENCY ( $\mathcal{C}$ )

INPUT: A DTD  $D$ , a set  $\Sigma$  of  $\mathcal{C}$ -constraints.

QUESTION: Is there an XML document that conforms to  $D$  and satisfies  $\Sigma$ ?

In other words, we want to validate XML specifications statically. The main reason is twofold:

---

<sup>1</sup>Throughout the paper, by a DTD we mean its type specification; we ignore its ID/IDREF constraints since their limitations have been well recognized [7; 19]. We shall only consider *finite* XML documents (trees).

first, complex interactions between DTDs and constraints are likely to result in inconsistent specifications, and second, an alternative dynamic approach to validation (simply check a document to see if it conforms to the DTD and satisfies the constraints) would not tell us whether repeated failures are due to a bad specification, or problems with the documents.

The concept of consistency of specifications was studied for other data models, such as object-oriented [12; 13] and extended relational (e.g., with support for cardinality constraints [26]).

We shall study the following four classes of constraints defined in terms of XML attributes:

- $\mathcal{C}_{K,FK}$ : a class of keys and foreign keys;
- $\mathcal{C}_{K,FK}^{Unary}$ : unary keys and foreign keys in  $\mathcal{C}_{K,FK}$ , i.e., those defined in terms of a single attribute;
- $\mathcal{C}_{K^-,IC}^{Unary}$ : unary keys, unary inclusion constraints and negations of unary keys;
- $\mathcal{C}_{K^-,IC^-}^{Unary}$ : unary keys, unary inclusion constraints and their negations.

Keys and foreign keys of  $\mathcal{C}_{K,FK}$  are a natural generalization of their relational counterpart, and are capable of capturing those relational constraints. A foreign key is a combination of two constraints: an inclusion constraint and a key. The  $\mathcal{C}_{K,FK}^{Unary}$  constraints are a special case of  $\mathcal{C}_{K,FK}$  constraints, which involve a single attribute. These unary keys and foreign keys are similar to but more general than XML ID and IDREF specifications. The study on simple constraints defined with XML attributes is a first step towards understanding the interaction between integrity constraints and schema specifications for XML. As will be seen shortly, the analyses of these simple constraints are already very intricate in the presence of DTDs.

As generalizations of  $\mathcal{C}_{K,FK}^{Unary}$  constraints,  $\mathcal{C}_{K^-,IC}^{Unary}$  and  $\mathcal{C}_{K^-,IC^-}^{Unary}$  both allow the presence of unary inclusion constraints independent of keys. In addition,  $\mathcal{C}_{K^-,IC}^{Unary}$  includes negations of unary keys, and  $\mathcal{C}_{K^-,IC^-}^{Unary}$  further permits negations of unary inclusion constraints. Negation is considered mainly for the study of *implication* of  $\mathcal{C}_{K,FK}^{Unary}$  constraints, which is the complement of a special case of the consistency problem for  $\mathcal{C}_{K^-,IC}^{Unary}$  (resp.  $\mathcal{C}_{K^-,IC^-}^{Unary}$ ): given any DTD  $D$  and any finite set  $\Sigma$  of unary keys and inclusion constraints, is it the case that all XML trees satisfying  $\Sigma$  and conforming to  $D$  also satisfy some other unary key (resp. unary key or inclusion constraint)? This question is important in, among other things, data integration. For example, one may want to know whether a constraint  $\varphi$  holds in a mediator interface, which may use XML as a uniform data format [4; 33]. This cannot be verified directly since the mediator interface does not contain data. One way to verify  $\varphi$  is to show that it is implied by constraints that are known to hold [21].

These problems, however, turn out to be far more intriguing than their counterparts in relational databases. In the XML setting, DTDs do interact with keys and foreign keys, and this interaction may lead to problems with XML specifications.

**Examples.** To illustrate the interaction between XML DTDs and key/foreign key constraints, consider a DTD  $D_1$ , which specifies a (nonempty) collection of teachers:

```
<!ELEMENT  teachers  (teacher+)>
<!ELEMENT  teacher   (teach, research)>
<!ELEMENT  teach     (subject, subject)>
```

It says that a teacher teaches two subjects. Here we omit the descriptions of elements whose type is string (e.g., PCDATA in XML).

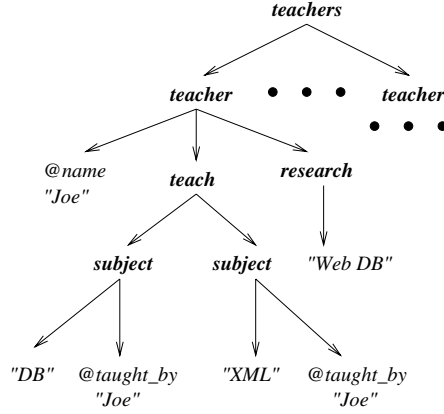
Assume that each teacher has an attribute **name** and each subject has an attribute **taught\_by**. Attributes are single-valued. That is, if an attribute  $l$  is defined for an element type  $\tau$  in a DTD, then in a document conforming to the DTD, each element of type  $\tau$  must have a unique  $l$  attribute with a string value. Consider a set of unary key and foreign key constraints,  $\Sigma_1$ :

$$\begin{aligned} & \text{teacher.name} \rightarrow \text{teacher}, \\ & \text{subject.taught\_by} \rightarrow \text{subject}, \\ & \text{subject.taught\_by} \subseteq \text{teacher.name}. \end{aligned}$$

That is, **name** is a key of **teacher** elements, **taught\_by** is a key of **subject** elements and it is also a foreign key referencing **name** of **teacher** elements. More specifically, referring to an XML tree  $T$ , the first constraint asserts that two distinct **teacher** nodes in  $T$  cannot have the same **name** attribute value: the (string) value of **name** attribute uniquely identifies a **teacher** node. It should be mentioned that two notions of equality are used in the definition of keys: we assume string *value* equality when comparing **name** attribute values, and *node* identity when it comes to comparing **teacher** elements. The second key states that **taught\_by** attribute uniquely identifies a **subject** node in  $T$ . The third constraint asserts that for any **subject** node  $x$ , there is a **teacher** node  $y$  in  $T$  such that the **taught\_by** attribute value of  $x$  equals the **name** attribute value of  $y$ . Since **name** is a key of **teacher**, the **taught\_by** attribute of any **subject** node refers to a unique **teacher** node.

Obviously, there exists an XML tree conforming to  $D_1$ , as shown in Figure 1. However, there is no XML tree that both conforms to  $D_1$  and satisfies  $\Sigma_1$ . To see this, let us first define some notations. Given an XML tree  $T$  and an element type  $\tau$ , we use  $\text{ext}(\tau)$  to denote the set of all the nodes labeled  $\tau$  in  $T$ . Similarly, given an attribute  $l$  of  $\tau$ , we use  $\text{ext}(\tau.l)$  to denote the set of  $l$  attribute values of all  $\tau$  elements. Then immediately from  $\Sigma_1$  follows a set of dependencies:

$$\begin{aligned} |\text{ext}(\text{teacher.name})| &= |\text{ext}(\text{teacher})|, \\ |\text{ext}(\text{subject.taught\_by})| &= |\text{ext}(\text{subject})|, \\ |\text{ext}(\text{subject.taught\_by})| &\leq |\text{ext}(\text{teacher.name})|, \end{aligned}$$

Fig. 1. An XML tree conforming to  $D_1$ 

where  $|\cdot|$  is the cardinality of a set. Therefore, we have

$$|ext(subject)| \leq |ext(teacher)|. \quad (1)$$

On the other hand, the DTD  $D_1$  requires that each teacher must teach two subjects. Since no sharing of nodes is allowed in XML trees and the collection of **teacher** elements is nonempty, from  $D_1$  follows:

$$1 < 2 |ext(teacher)| = |ext(subject)|. \quad (2)$$

Thus  $|ext(teacher)| < |ext(subject)|$ . Obviously, (1) and (2) contradict with each other and therefore, there exists no XML tree that both satisfies  $\Sigma_1$  and conforms to  $D_1$ . In particular, the XML tree in Figure 1 violates the key  $subject.taught\_by \rightarrow subject$ .

This example demonstrates that a DTD may impose dependencies on the cardinalities of certain sets of objects in XML trees. These *cardinality constraints* interact with keys and foreign keys. More specifically, keys and foreign keys also enforce cardinality constraints that interact with those imposed by DTD. This makes the consistency analysis of keys and foreign keys for XML far more intriguing than that for relational databases. Because of the interaction, simple key and foreign key constraints (e.g.,  $\Sigma_1$ ) may not be satisfiable by XML trees conforming to certain DTDs (e.g.,  $D_1$ ).

As another example, consider the DTD  $D_2$  given below:

```
<!ELEMENT db (foo)>
<!ELEMENT foo (foo)>
```

Observe that there exists no finite XML tree conforming to  $D_2$ . This demonstrates that there is need for studying consistency of XML specifications even in the absence of integrity constraints.

**Contributions.** The main contributions of the paper are the following:

- (1) For the class  $\mathcal{C}_{K,FK}$  of keys and foreign keys, we show that both the consistency and the implication problems are undecidable.
- (2) These negative results suggest that we look at the restriction  $\mathcal{C}_{K,FK}^{Unary}$  of *unary* keys and foreign keys (which are most typical in XML documents). We provide a coding of DTDs and these unary constraints by linear constraints on the integers. This enables us to show that the consistency problem for  $\mathcal{C}_{K,FK}^{Unary}$  (even under the restriction to primary keys, i.e., at most one key for each element type) is NP-complete. We further show that the problem is still in NP for an extension  $\mathcal{C}_{K^-,IC}^{Unary}$ , which also allows negations of key constraints.
- (3) Using a different coding of constraints, we show that the consistency problem remains in NP for  $\mathcal{C}_{K^-,IC^-,}^{Unary}$ , the class of unary keys, unary inclusion constraints and their negations. Among other things, this shows that the implication problem for unary keys and foreign keys is coNP-complete.
- (4) We also identify several tractable cases of the consistency problem, i.e., practical situations where the consistency problem is decidable in PTIME.

The undecidability of the consistency problem contrasts sharply with its trivial counterpart in relational databases. The coding of DTDs and unary constraints with linear integer constraints reveals some insight into the interaction between DTDs and unary constraints. Moreover, it allows us to use the techniques from linear integer programming in the study of XML constraints.

It should be mentioned that as XML Schema and XML Data both subsume DTDs and they support keys and foreign keys which are more general than those considered here, the undecidability and NP-hardness results carry over to these schema specifications and constraint languages for XML.

**Related work.** Keys, foreign keys and the more general inclusion and functional dependencies have been well studied for relational databases (cf. [1]). In particular, the implication problem for unary inclusion and functional dependencies is in linear time [17]. In contrast, we shall show that the XML counterpart of this problem is coNP-complete.

The interaction between cardinality constraints and database schemas has been studied for object-oriented [12; 13] and extended relational data models [26]. These interactions are quite different from what we explore in this paper because XML DTDs are defined in terms of extended context free grammars and they yield cardinality constraints more complex than those studied for databases.

Key and foreign key specifications for XML have been proposed in the XML standard [6], XML Data [27], XML Schema [36] and in a recent proposal for XML keys [7]. The need for

studying XML constraints has also been advocated in [38]. DTDs in the XML standard allow one to specify limited (primary) unary keys and foreign keys with ID and IDREF attributes. However, they are not scoped: one has no control over what IDREF attributes point to. XML Data and XML Schema support more expressive specifications for keys and foreign keys with, e.g., XPath expressions. However, the consistency problems associated with constraints defined in these languages have not been studied. We consider simple XML keys and foreign keys in this paper to focus on the nature of the interaction between DTDs and constraints. The implication problem for a class of keys and foreign keys was investigated in [19], but in the absence of DTDs (in a graph model for XML), which trivializes the consistency analysis. For keys of [7], the implication problem was studied [8] in the tree model for XML, but DTDs were not considered there. To the best of our knowledge, no previous work has considered the interaction between DTDs and keys and foreign keys for XML (in the tree model). This paper is a full version of [18], providing the details and the proofs omitted there.

A variety of constraints have been studied for semistructured data [2; 10; 20]. In particular, [20] also studies the consistency problem; the special form of constraints used there makes it possible to encode consistency as an instance of conjunctive query containment. The interaction between path constraints and database schemas was investigated in [9]. These constraints typically specify inclusions among certain sets of objects in edge-labeled graphs, and are not capable of expressing keys. Various generalizations of functional dependencies have also been studied [23; 25]. But these generalizations were investigated in database settings, which are quite different from the tree model for XML data. Moreover, they cannot express foreign keys. Application of constraints in data transformations was studied in [28]; usefulness of keys and foreign keys in query optimization has also been recognized [34].

**Organization.** The rest of the paper is organized as follows. Section 2 defines four classes of XML constraints, namely,  $\mathcal{C}_{K,FK}$ ,  $\mathcal{C}_{K,FK}^{Unary}$ ,  $\mathcal{C}_{K^-,IC}^{Unary}$  and  $\mathcal{C}_{K^-,IC^-}^{Unary}$ . Section 3 establishes the undecidability of the consistency problem for  $\mathcal{C}_{K,FK}$ , the class of keys and foreign keys. Section 4 provides an encoding for DTDs and unary constraints with linear integer constraints, and shows that the consistency problems are NP-complete for  $\mathcal{C}_{K,FK}^{Unary}$  and  $\mathcal{C}_{K^-,IC}^{Unary}$ . Section 5 further shows that the problem remains in NP for  $\mathcal{C}_{K^-,IC^-}^{Unary}$ , the class of unary keys, inclusion constraints and their negations. Section 6 summarizes the main results of the paper and identifies directions for further work.

## 2. DTDs, KEYS AND FOREIGN KEYS

In this section, we first present a formalism of XML DTDs [6] and the XML tree model. We then define four classes of XML constraints.

## 2.1 DTDs and XML trees

We extend the usual formalism of DTDs (as extended context free grammars [5; 11; 31]) by incorporating attributes.

**DEFINITION 2.1.** *A DTD (Document Type Definition) is defined to be  $D = (E, A, P, R, r)$ , where:*

- $E$  is a finite set of element types;
- $A$  is a finite set of attributes, disjoint from  $E$ ;
- $P$  is a mapping from  $E$  to element type definitions: for each  $\tau \in E$ ,  $P(\tau)$  is a regular expression  $\alpha$  defined as follows:

$$\alpha ::= \mathbf{S} \mid \tau' \mid \epsilon \mid \alpha|\alpha \mid \alpha, \alpha \mid \alpha^*$$

where  $\mathbf{S}$  denotes string type,  $\tau' \in E$ ,  $\epsilon$  is the empty word, and “|”, “,” and “\*” denote union, concatenation, and the Kleene closure, respectively;

- $R$  is a mapping from  $E$  to  $\mathcal{P}(A)$ , the power-set of  $A$ ; if  $l \in R(\tau)$  then we say  $l$  is defined for  $\tau$ ;
- $r \in E$  and is called the element type of the root.

We normally denote element types by  $\tau$  and attributes by  $l$ . Without loss of generality, assume that  $r$  does not occur in  $P(\tau)$  for any  $\tau \in E$ . We also assume that each  $\tau$  in  $E \setminus \{r\}$  is connected to  $r$ , i.e., either  $\tau$  occurs in  $P(r)$ , or it appears in  $P(\tau')$  for some  $\tau'$  that is connected to  $r$ .

As an example, let us consider the **teacher** DTD  $D_1$  given in Section 1. In our formalism,  $D_1$  can be represented as  $(E_1, A_1, P_1, R_1, r_1)$ , where

$$\begin{aligned} E_1 &= \{\text{teachers}, \text{teacher}, \text{teach}, \text{research}, \text{subject}\} \\ A_1 &= \{\text{name}, \text{taught\_by}\} \\ P_1(\text{teachers}) &= \text{teacher}, \text{teacher}^* \\ P_1(\text{teacher}) &= \text{teach}, \text{research} \\ P_1(\text{teach}) &= \text{subject}, \text{subject} \\ P_1(\text{subject}) &= P_1(\text{research}) = \mathbf{S} \\ R_1(\text{teacher}) &= \{\text{name}\} \\ R_1(\text{subject}) &= \{\text{taught\_by}\} \\ R_1(\text{teachers}) &= R_1(\text{teach}) = R_1(\text{research}) = \emptyset \\ r_1 &= \text{teachers} \end{aligned}$$

Similarly, we represent the DTD  $D_2$  given in Section 1 as  $(E_2, A_2, P_2, R_2, r_2)$ , where



$$\begin{aligned}
E_2 &= \{db, foo\} \\
A_2 &= \emptyset \\
P_2(db) &= P_2(foo) = foo \\
R_2(db) &= R_2(foo) = \emptyset \\
r_2 &= db
\end{aligned}$$

An XML document is typically modeled as a node-labeled ordered tree. Given a DTD, we define the notion of its valid documents as follows.

**DEFINITION 2.2.** *Let  $D = (E, A, P, R, r)$  be a DTD. An XML tree  $T$  valid w.r.t.  $D$  (conforming to  $D$ ) is defined to be  $T = (V, lab, ele, att, val, root)$ , where*

- $V$  is a finite set of nodes (vertices);
- $lab$  is a function that maps each node in  $V$  to a label in  $E \cup A \cup \{\mathbf{S}\}$ ; a node  $v \in V$  is called an element of  $\tau$  if  $lab(v) = \tau$  and  $\tau \in E$ , an attribute if  $lab(v) \in A$ , and a text node if  $lab(v) = \mathbf{S}$ ;
- $ele$  is a partial function defined on elements in  $V$ ; for any  $\tau \in E$ , it maps each element  $v$  of type  $\tau$  to a (possibly empty) list  $[v_1, \dots, v_n]$  of elements and text nodes in  $V$  such that  $lab(v_1) \dots lab(v_n)$  is in the regular language defined by  $P(\tau)$ ;
- $att$  is a partial function from  $V \times A$  to  $V$  such that for any  $v \in V$  and  $l \in A$ ,  $att(v, l)$  is defined iff  $lab(v) = \tau$ ,  $\tau \in E$  and  $l \in R(\tau)$ ;
- $val$  is a partial function from  $V$  to string values such that for any node  $v \in V$ ,  $val(v)$  is defined iff  $lab(v) = \mathbf{S}$  or  $lab(v) \in A$ ;
- $root$  is the unique node in  $V$  such that  $lab(root) = r$ , called the root of  $T$ .

For any element  $v \in V$ , the nodes  $v'$  in  $ele(v)$  are called the subelements of  $v$ . For any  $l \in A$ , if  $att(v, l) = v'$  then  $v'$  is called an attribute of  $v$ . In either case we say that there is a parent-child edge from  $v$  to  $v'$ . The subelements and attributes of  $v$  are called its children. An XML tree has a tree structure, i.e., for each  $v \in V$ , there is a unique path of parent-child edges from root to  $v$ . We write  $T \models D$  when  $T$  is valid w.r.t.  $D$ .

Intuitively,  $V$  is the set of nodes of the tree  $T$ . The mapping  $lab$  labels every node of  $V$  with a symbol from  $E \cup A \cup \{\mathbf{S}\}$ . Text nodes and attributes are leaves. For an element  $x$  of type  $\tau$ , the functions  $ele$  and  $att$  define the children of  $x$ , which are partitioned into *subelements* and *attributes* according to  $P(\tau)$  and  $R(\tau)$  in the DTD  $D$ . The subelements of  $x$  are ordered and their labels satisfy the regular expression  $P(\tau)$ . In contrast, its attributes are unordered and are identified by their labels (names). The function  $val$  assigns string values to attributes and text nodes. We consider single-valued attributes. That is, if  $l \in R(\tau)$  then every element of type  $\tau$  has a unique  $l$  attribute with a string value. Since  $T$  has a tree structure, sharing of nodes is not allowed in  $T$ .

For example, Figure 1 depicts an XML tree valid w.r.t. the DTD  $D_1$  given in Section 1.

Our model is simpler than the models of XQuery [14] and XML Schema [36] as DTDs support only one basic type (PCDATA or string) and do not have complex type constructs. Furthermore, we do not have nodes representing namespaces, processing instructions and references. These simplifications allow us to concentrate on the essence of the DTD/constraint interaction. It should further be noticed that they do not affect the lower bounds results in the paper.

We need the following notations throughout the paper: for any  $\tau \in E \cup \{\mathbf{S}\}$ ,  $ext(\tau)$  denotes the set of all the nodes in  $T$  labeled  $\tau$ . For any node  $x$  in  $T$  labeled by  $\tau$  and for any attribute  $l \in R(\tau)$ , we write  $x.l$  for  $val(att(x, l))$ , i.e., the *value* of the attribute  $l$  of node  $x$ . We define  $ext(\tau.l)$  to be  $\{x.l \mid x \in ext(\tau)\}$ , which is a set of strings. For each  $\tau$  element  $x$  in  $T$  and a list  $X = [l_1, \dots, l_n]$  of attributes in  $R(\tau)$ , we use  $x[X]$  to denote the list of  $X$ -attribute values of  $x$ , i.e.,  $x[X] = [x.l_1, \dots, x.l_n]$ . For a set  $S$ ,  $|S|$  denotes its cardinality.

## 2.2 XML constraints

We next define our constraint languages for XML.

We consider three types of constraints. Let  $D = (E, A, P, R, r)$  be a DTD, and  $T$  be an XML tree valid w.r.t.  $D$ . A *constraint*  $\varphi$  over  $D$  has one of the following forms:

—*Key*:  $\tau[X] \rightarrow \tau$ , where  $\tau \in E$  and  $X$  is a set of attributes in  $R(\tau)$ . The XML tree  $T$  *satisfies*  $\varphi$ , denoted by  $T \models \varphi$ , iff in  $T$ ,

$$\forall x y \in ext(\tau) \left( \bigwedge_{l \in X} (x.l = y.l) \rightarrow x = y \right).$$

—*Inclusion constraint*:  $\tau_1[X] \subseteq \tau_2[Y]$ , where  $\tau_1, \tau_2 \in E$ , and  $X, Y$  are nonempty lists of attributes in  $R(\tau_1), R(\tau_2)$  of the same length. We write  $T \models \varphi$  iff in  $T$ ,

$$\forall x \in ext(\tau_1) \exists y \in ext(\tau_2) (x[X] = y[Y]).$$

—*Foreign key*: a combination of two constraints, namely, an inclusion constraint  $\tau_1[X] \subseteq \tau_2[Y]$  and a key  $\tau_2[Y] \rightarrow \tau_2$ . We write  $T \models \varphi$  iff  $T$  satisfies both the key and the inclusion constraint.

That is, a key  $\tau[X] \rightarrow \tau$  indicates that the set  $X$  of attributes is a key of elements of  $\tau$ , i.e., two distinct  $\tau$  nodes in  $T$  cannot have the same  $X$ -attribute values; an inclusion constraint  $\tau_1[X] \subseteq \tau_2[Y]$  says that the list of  $X$ -attribute values of every  $\tau_1$  node in  $T$  must match the list of  $Y$ -attribute values of some  $\tau_2$  node in  $T$ ; and an foreign key  $\tau_1[X] \subseteq \tau_2[Y]$ ,  $\tau_2[Y] \rightarrow \tau_2$  indicates that  $X$  is a foreign key of  $\tau_1$  elements referencing key  $Y$  of  $\tau_2$  elements.

Over a DTD  $D$ , the class  $\mathcal{C}_{K,FK}$  of constraints consists of all the keys and foreign keys over  $D$ . They are called *multi-attribute* keys and foreign keys as they may be defined in terms of multiple attributes.

To illustrate keys and foreign keys of  $\mathcal{C}_{K,FK}$ , let us consider a DTD  $D_3 = (E_3, A_3, P_3, R_3, r_3)$ , where

$$\begin{aligned}
E_3 &= \{school, student, course, enroll, name, subject\} \\
A_3 &= \{student\_id, course\_no, dept\} \\
P_3(school) &= course^*, student^*, enroll^* \\
P_3(course) &= subject \\
P_3(student) &= name \\
P_3(enroll) &= P_3(name) = P_3(subject) = S \\
R_3(course) &= \{dept, course\_no\} \\
R_3(student) &= \{student\_id\} \\
R_3(enroll) &= \{student\_id, dept, course\_no\} \\
R_3(school) &= R_3(name) = R_3(subject) = \emptyset \\
r_3 &= school
\end{aligned}$$

Typical  $\mathcal{C}_{K,FK}$  constraints over  $D_3$  include:

- (1)  $student[student\_id] \rightarrow student$ ,
- (2)  $course[dept, course\_no] \rightarrow course$ ,
- (3)  $enroll[student\_id, dept, course\_no] \rightarrow enroll$ ,
- (4)  $enroll[student\_id] \subseteq student[student\_id]$ ,
- (5)  $enroll[dept, course\_no] \subseteq course[dept, course\_no]$ .

The first three constraints are keys in  $\mathcal{C}_{K,FK}$ , and the pairs (4, 1) and (5, 2) are foreign keys in  $\mathcal{C}_{K,FK}$ . The last two constraints are inclusion constraints.

It is worth mentioning that two notions of equality are used to define keys: string value equality is assumed in  $x.l = y.l$  (when comparing attribute values), and  $x = y$  is true if and only if  $x$  and  $y$  are the same node (when comparing elements). This is different from the semantics of keys in relational databases. Note that a foreign key requires the presence of a key in addition to an inclusion constraint.

The class of unary keys and foreign keys for XML, denoted by  $\mathcal{C}_{K,FK}^{Unary}$ , is a sublanguage of  $\mathcal{C}_{K,FK}$ . A  $\mathcal{C}_{K,FK}^{Unary}$  constraint is a  $\mathcal{C}_{K,FK}$  constraint defined with a single attribute. More specifically, a constraint  $\varphi$  of  $\mathcal{C}_{K,FK}^{Unary}$  over the DTD  $D$  is either

- key*:  $\tau.l \rightarrow \tau$ , where  $\tau \in E$  and  $l \in R(\tau)$ ; or
- foreign key*:  $\tau_1.l_1 \subseteq \tau_2.l_2$  and  $\tau_2.l_2 \rightarrow \tau_2$ , where  $\tau_1, \tau_2 \in E$ ,  $l_1 \in R(\tau_1)$ , and  $l_2 \in R(\tau_2)$ .

For example, the constraints of  $\Sigma_1$  given in Section 1 are  $\mathcal{C}_{K,FK}^{Unary}$  constraints over the DTD  $D_1$ .

We shall also consider the following types of unary constraints over  $D$ :

- inclusion constraint*:  $\tau_1.l_1 \subseteq \tau_2.l_2$ ; unlike a foreign key, it does not require the presence of a key;
- the *negation of an inclusion constraint*:  $\phi = \tau_1.l_1 \not\subseteq \tau_2.l_2$ ; for an XML tree  $T$ ,  $T \models \phi$  iff there is a  $\tau_1$  element  $x$  in  $T$  such that for all  $\tau_2$  element  $y$  in  $T$ ,  $x.l_1 \neq y.l_2$ ;
- the *negation of a key*:  $\varphi = \tau.l \not\rightarrow \tau$ ;  $T \models \varphi$  iff there are  $\tau$  elements  $x_1, x_2$  in  $T$  such that  $x_1.l = x_2.l$ , i.e., the value of the  $l$  attribute of a  $\tau$  element cannot uniquely identify it in  $ext(\tau)$ .

With these we define two extensions of  $\mathcal{C}_{K,FK}^{Unary}$  as follows. One is  $\mathcal{C}_{K^-,IC}^{Unary}$ , the class consisting of unary keys, unary inclusion constraints and negations of unary keys. The other,  $\mathcal{C}_{K^-,IC^-}^{Unary}$ , consists of unary keys, unary inclusion constraints and their negations. As mentioned earlier, we consider these classes mostly for the study of the implication problem for  $\mathcal{C}_{K,FK}^{Unary}$  constraints.

Finally, we describe the consistency and implication problems associated with XML constraints. Let  $\mathcal{C}$  be one of  $\mathcal{C}_{K,FK}$ ,  $\mathcal{C}_{K,FK}^{Unary}$ ,  $\mathcal{C}_{K^-,IC}^{Unary}$  or  $\mathcal{C}_{K^-,IC^-}^{Unary}$ ,  $D$  a DTD,  $\Sigma$  a set of  $\mathcal{C}$  constraints over  $D$  and  $T$  an XML tree valid w.r.t.  $D$ . We write  $T \models \Sigma$  when  $T \models \phi$  for all  $\phi \in \Sigma$ . Let  $\varphi$  be another  $\mathcal{C}$  constraint. We say that  $\Sigma$  *implies*  $\varphi$  over  $D$ , denoted by  $(D, \Sigma) \vdash \varphi$ , if for any XML tree  $T$  such that  $T \models D$  and  $T \models \Sigma$ , it must be the case that  $T \models \varphi$ . It should be noted when  $\varphi$  is a foreign key,  $\varphi$  consists of an inclusion constraint  $\phi_1$  and a key  $\phi_2$ . In this case  $(D, \Sigma) \vdash \varphi$  in fact means that  $(D, \Sigma) \vdash \phi_1 \wedge \phi_2$ .

The central technical problem investigated in this paper is the *consistency problem*. The consistency problem for  $\mathcal{C}$  is to determine, given any DTD  $D$  and any set  $\Sigma$  of  $\mathcal{C}$  constraints over  $D$ , whether there is an XML tree  $T$  such that  $T \models \Sigma$  and  $T \models D$ .

The *implication problem* for  $\mathcal{C}$  is to determine, given any DTD  $D$ , any set  $\Sigma$  and  $\varphi$  of  $\mathcal{C}$  constraints over  $D$ , whether  $(D, \Sigma) \vdash \varphi$ .

### 3. GENERAL KEYS AND FOREIGN KEYS

In this section we study  $\mathcal{C}_{K,FK}$ , the class of multi-attribute keys and foreign keys. We show that the consistency and implication problems for  $\mathcal{C}_{K,FK}$  are undecidable, but we identify several special cases of the problems and show that these cases are decidable in PTIME.

### 3.1 Undecidability of consistency analysis

Our main result is negative:

**THEOREM 3.1.** *The consistency problem for  $\mathcal{C}_{K,FK}$  constraints is undecidable.*

*Proof:* We first show that an implication problem associated with keys and foreign keys in relational databases is undecidable, and then present a reduction from (the complement of) the implication problem to the consistency problem for  $\mathcal{C}_{K,FK}$  constraints.

Let us first review keys, foreign keys and their associated implication problems in relational databases (cf. [1]). Let  $\mathbf{R} = (R_1, \dots, R_n)$  be a relational schema. For each relation (schema)  $R_i$  in  $\mathbf{R}$ , we write  $Att(R_i)$  for the set of all attributes of  $R_i$ , and  $Inst(R_i)$  for the set of finite instances of  $R_i$ . By database instances we mean *finite* instances. An instance  $\mathbf{I}$  of  $\mathbf{R}$  has the form  $(I_1, \dots, I_n)$ , where  $I_i \in Inst(R_i)$  for all  $i \in [1, n]$ . For an instance  $I_i \in Inst(R_i)$ , a tuple  $t \in I_i$  and an attribute  $l \in Att(R_i)$ , we use  $t.l$  to denote the  $l$  attribute value of  $t$ . Keys and foreign keys over  $\mathbf{R}$  are defined as follows:

—*key:*  $R[l_1, \dots, l_k] \rightarrow R$ , where  $R \in \mathbf{R}$ , and for any  $i \in [1, k]$ ,  $l_i \in Att(R)$ . An instance  $\mathbf{I}$  of  $\mathbf{R}$  *satisfies* the key constraint  $\varphi$ , denoted by  $\mathbf{I} \models \varphi$ , if

$$\forall t_1 t_2 \in I \left( \bigwedge_{1 \leq i \leq k} (t_1.l_i = t_2.l_i) \rightarrow \bigwedge_{l \in Att(R)} (t_1.l = t_2.l) \right),$$

where  $I$  is the instance of  $R$  in  $\mathbf{I}$ ;

—*foreign key:*  $R[l_1, \dots, l_k] \subseteq R'[l'_1, \dots, l'_k]$  and  $R'[l'_1, \dots, l'_k] \rightarrow R'$ , where  $R, R'$  are in  $\mathbf{R}$ ,  $[l_1, \dots, l_k]$  and  $[l'_1, \dots, l'_k]$  are lists of attributes in  $Att(R)$  and in  $Att(R')$ , respectively. In addition, the set consisting of  $l'_1, \dots, l'_k$  is a key of  $R'$ . We write  $\mathbf{I} \models \varphi$  if  $\mathbf{I} \models R'[l'_1, \dots, l'_k] \rightarrow R'$  and moreover,

$$\forall t_1 \in I \exists t_2 \in I' \left( \bigwedge_{1 \leq j \leq k} t_1.l_j = t_2.l'_j \right),$$

where  $I$  and  $I'$  are the instances of  $R$  and  $R'$  in  $\mathbf{I}$ , respectively.

Let  $\Sigma \cup \{\varphi\}$  be a set of keys and foreign keys over  $\mathbf{R}$ . We use  $\Sigma \vdash \varphi$  to denote that  $\Sigma$  *implies*  $\varphi$ , i.e., for any instance  $\mathbf{I}$  of  $\mathbf{R}$ , if  $\mathbf{I} \models \Sigma$ , then  $\mathbf{I} \models \varphi$ .

In relational databases, the *implication problem for keys and foreign keys* is the problem of determining, given a relational schema  $\mathbf{R}$ , any set  $\Sigma$  and  $\varphi$  of keys and foreign keys over  $\mathbf{R}$ , whether  $\Sigma \vdash \varphi$ . A special case of the problem is the *implication problem for keys by keys and foreign keys*, which is to determine whether  $\Sigma \vdash \varphi$  where  $\varphi$  is a key and  $\Sigma$  is a set of keys and foreign keys over  $\mathbf{R}$ .

It was shown in [19] that the implication problem for keys and foreign keys in relational databases is undecidable. The lemma below shows a stronger result.

**LEMMA 3.2.** *In relational databases, the implication problem for keys by keys and foreign keys is undecidable.*

*Proof:* We prove this by reduction from the implication problem for functional dependencies by functional and inclusion dependencies, which is undecidable. Before we give the reduction, we first review functional and inclusion dependencies in relational databases. Let  $\mathbf{R}$  be a relational schema. Functional dependencies (FDs) and inclusion dependencies (IDs) over  $\mathbf{R}$  are defined as follows.

- FD.*  $R : X \rightarrow Y$ , where  $R \in \mathbf{R}$ , and  $X$  and  $Y$  are subsets of attributes in  $Att(R)$ . An instance  $\mathbf{I}$  of  $\mathbf{R}$  *satisfies* the FD  $\theta$ , denoted by  $\mathbf{I} \models \theta$ , if  $\forall t_1 t_2 \in I \left( \bigwedge_{l \in X} (t_1.l = t_2.l) \rightarrow \bigwedge_{l' \in Y} (t_1.l' = t_2.l') \right)$ , where  $I$  is the instance of  $R$  in  $\mathbf{I}$ . Observe that keys are a special case of FDs in which  $Y = Att(R)$ .
- ID.*  $R[l_1, \dots, l_k] \subseteq R'[l'_1, \dots, l'_k]$ , where  $R, R' \in \mathbf{R}$ ,  $[l_1, \dots, l_k]$  is a list of attributes in  $Att(R)$ , and  $[l'_1, \dots, l'_k]$  is a list of attributes in  $Att(R')$ . In contrast to foreign keys, the set consisting of  $l'_1, \dots, l'_k$  is not necessarily a key of  $R'$ . An instance  $\mathbf{I}$  of  $\mathbf{R}$  *satisfies* the ID  $\theta$ , denoted by  $\mathbf{I} \models \theta$ , if  $\forall t_1 \in I \exists t_2 \in I' \left( \bigwedge_{1 \leq j \leq k} t_1.l_j = t_2.l'_j \right)$ , where  $I, I'$  are the instances of  $R, R'$  in  $\mathbf{I}$ , respectively.

Let  $\Sigma \cup \{\theta\}$  be a set of FDs and IDs over  $\mathbf{R}$ . We use  $\Sigma \vdash \theta$  to denote that  $\Sigma$  *implies*  $\theta$  as for keys and foreign keys. The *implication problem for FDs by FDs and IDs* is the problem to determine, given any relational schema  $\mathbf{R}$ , any set  $\Sigma$  of FDs and IDs over  $\mathbf{R}$  and a FD  $\theta$  over  $\mathbf{R}$ , whether  $\Sigma \vdash \theta$ . This is a well-known undecidable problem (see, e.g., [1] for a proof).

We encode FDs and IDs in terms of keys and foreign keys as follows.

- (1) FD  $\psi = R : X \rightarrow Y$ .

Note that every relation  $R$  has a key. In particular,  $Att(R)$ , the set of all attributes of  $R$ , is a key of  $R$ . Let  $Z$  be a key for  $R$ , i.e.,  $R[Z] \rightarrow R$ . We define a new (fresh) relation schema  $R_{new}$  such that  $Att(R_{new}) = XYZ$ , i.e., the union of  $X, Y$  and  $Z$ . Intuitively, given an instance  $I$  of  $R$ , an instance  $I_{new}$  of  $R_{new}$  is to be constructed as a subset of  $\Pi_{XYZ}(I)$  such that  $\Pi_{XY}(I) = \Pi_{XY}(I_{new})$  and  $I_{new}$  satisfies the key  $R_{new}[XY] \rightarrow R_{new}$ , where  $\Pi_W(I)$  denotes the projection of  $I$  on attributes  $W$ . That is, we eliminate tuples in  $\Pi_{XYZ}(I)$  that violate the key. Observe that  $XYZ$  is a key for both  $R_{new}$  and  $R$  since it is the set of all attributes of  $R_{new}$ ,

and it contains the key  $Z$  of  $R$  (i.e., it is a *superkey* of  $R$ ). Thus we encode  $\psi$  with:

$$\begin{aligned}\phi_1 &= R_{new}[X] \rightarrow R_{new}, & \phi_2 &= R[XY] \subseteq R_{new}[XY], \\ \phi_3 &= R_{new}[XYZ] \subseteq R[XYZ], & \phi_4 &= R_{new}[XY] \rightarrow R_{new}.\end{aligned}$$

(2) ID  $\psi = R_1[X] \subseteq R_2[Y]$ .

Let  $Z$  be a key for  $R_2$ , i.e.,  $R_2[Z] \rightarrow R_2$ . We define a new schema  $R_{new}$  such that  $Att(R_{new}) = YZ$ . Intuitively, given an instance  $I_2$  of  $R_2$ , an instance  $I_{new}$  of  $R_{new}$  is to be constructed as a subset of  $\Pi_{YZ}(I_2)$  by eliminating tuples that violate the key  $R_{new}[Y] \rightarrow R_{new}$ , such that  $\Pi_Y(I_2) = \Pi_Y(I_{new})$  and  $I_{new}$  satisfies the key. Observe that  $YZ$  is a key for  $R_2$  since it contains the key  $Z$  of  $R_2$ , i.e., it is a superkey of  $R_2$ . Thus we encode  $\psi$  with:

$$\phi_1 = R_{new}[Y] \rightarrow R_{new}, \quad \phi_2 = R_1[X] \subseteq R_{new}[Y], \quad \phi_3 = R_{new}[YZ] \subseteq R_2[YZ].$$

We next show that the encoding is indeed a reduction from the implication problem for FDs by FDs and IDs to the implication problem for keys by keys and foreign keys. Given a relational schema  $\mathbf{R}$ , a set  $\Sigma$  of FDs and IDs over  $\mathbf{R}$ , and a FD  $\theta = R_\theta : X \rightarrow Y$  over  $\mathbf{R}$ , as described above we encode  $\Sigma$  with a set  $\Sigma_1$  of keys and foreign keys, and encode  $\theta$  with

$$\begin{aligned}\phi_1 &= R_{new}^\theta[X] \rightarrow R_{new}^\theta, & \phi_2 &= R_\theta[XY] \subseteq R_{new}^\theta[XY], \\ \phi_3 &= R_{new}^\theta[XYZ] \subseteq R_\theta[XYZ], & \phi_4 &= R_{new}^\theta[XY] \rightarrow R_{new}^\theta.\end{aligned}$$

Let  $\Sigma' = \Sigma_1 \cup \{\phi_2, \phi_3, \phi_4\}$ . It suffices to show that  $\Sigma \vdash \theta$  iff  $\Sigma' \vdash \phi_1$ .

Let  $\mathbf{R}'$  be the relational schema that includes all relation schemas in  $\mathbf{R}$  as well as new relations created in the encoding. We show the claim as follows.

(1) Suppose that there is an instance  $\mathbf{I}$  of  $\mathbf{R}$  such that  $\mathbf{I} \models \bigwedge \Sigma \wedge \neg \theta$ . We show that there is an instance  $\mathbf{I}'$  of  $\mathbf{R}'$  such that  $\mathbf{I}' \models \bigwedge \Sigma' \wedge \neg \phi_1$ . We construct  $\mathbf{I}'$  such that for any  $R$  in  $\mathbf{R}$ , the instance of  $R$  in  $\mathbf{I}'$  is the same as the instance of  $R$  in  $\mathbf{I}$ . We populate instances of new relations  $R_{new}$  created in the encoding as mentioned above. (a) If  $R_{new}$  is introduced in the encoding of a FD  $R : X \rightarrow Y$  then we let the instance  $I_{new}$  of  $R_{new}$  in  $\mathbf{I}'$  be a subset of  $\Pi_{XYZ}(I)$  such that  $\Pi_{XY}(I) = \Pi_{XY}(I_{new})$  and  $I_{new} \models R_{new}[XY] \rightarrow R_{new}$ , where  $I$  is the instance of  $R$  in  $\mathbf{I}$ . (b) If  $R_{new}$  is introduced in the encoding of an ID  $R_1[X] \subseteq R_2[Y]$  then let the instance  $I_{new}$  of  $R_{new}$  in  $\mathbf{I}'$  be a subset of  $\Pi_{YZ}(I_2)$  such that  $\Pi_Y(I_2) = \Pi_Y(I_{new})$  and  $I_{new} \models R_{new}[Y] \rightarrow R_{new}$ , where  $I_2$  is the instance of  $R_2$  in  $\mathbf{I}$ . It is easy to verify that  $\mathbf{I}' \models \bigwedge \Sigma' \wedge \neg \phi_1$ .

(2) Suppose that there is an instance  $\mathbf{I}'$  of  $\mathbf{R}'$  such that  $\mathbf{I}' \models \bigwedge \Sigma' \wedge \neg \phi_1$ . We construct an instance  $\mathbf{I}$  of  $\mathbf{R}$  by removing from  $\mathbf{I}'$  all instances of new relations introduced in the encoding. It is easy to verify that  $\mathbf{I} \models \bigwedge \Sigma \wedge \neg \theta$ .

Therefore, the encoding is indeed a reduction from the implication problem for FDs by FDs and IDs. This shows that the implication problem for keys by keys and foreign keys is undecidable.  $\square$

From Lemma 3.2 follows that the complement of the implication problem for keys by keys and foreign keys is also undecidable. That is to determine, given a relational schema  $\mathbf{R}$ , a set  $\Sigma$  of keys and foreign keys over  $\mathbf{R}$  and a key  $\varphi$  over  $\mathbf{R}$ , whether there is an instance of  $\mathbf{R}$  satisfying  $\bigwedge \Sigma \wedge \neg \varphi$ .

We now continue with the proof of Theorem 3.1, i.e., the consistency problem for  $\mathcal{C}_{K,FK}$  constraints is undecidable. Given Lemma 3.2, it suffices to give a reduction from the complement of the implication problem for keys by keys and foreign keys. Let  $\mathbf{R} = (R_1, \dots, R_n)$  be a relational schema,  $\Theta$  be a set of keys and foreign keys over  $\mathbf{R}$ , and  $\varphi = R[X] \rightarrow R$  be a key over  $\mathbf{R}$ . Let  $Y = \text{Att}(R) \setminus X$ . We encode  $\mathbf{R}$ ,  $\Theta$  and  $\varphi$  in terms of a DTD  $D$  and a set  $\Sigma$  of  $\mathcal{C}_{K,FK}$  constraints over  $D$  as follows. Let  $D = (E, A, P, R_A, r)$ , where

$$\begin{aligned} E &= \{R_i \mid i \in [1, n]\} \cup \{t_i \mid i \in [1, n]\} \cup \{r, D_Y, E_X\} \\ A &= \bigcup_{i \in [1, n]} \text{Att}(R_i) \\ P(r) &= R_1, \dots, R_n, D_Y, D_Y, E_X \\ P(R_i) &= t_i^* \quad \text{for } i \in [1, n] \\ P(t_i) &= \epsilon \quad \text{for } i \in [1, n] \\ P(D_Y) &= P(E_X) = \epsilon \\ R_A(t_i) &= \text{Att}(R_i) \quad \text{for } i \in [1, n] \\ R_A(D_Y) &= X \cup Y \\ R_A(E_X) &= X \\ R_A(r) &= R_A(R_i) = \emptyset \quad \text{for } i \in [1, n] \end{aligned}$$

We denote  $P(R) = t_\varphi^*$  for the relation  $R$  in  $\varphi$ . Note that  $R = R_s$  and  $t_\varphi = t_s$  for some  $s \in [1, n]$ .

We encode  $\Theta$  and  $\varphi$  with  $\Sigma = \Sigma_\Theta \cup \Sigma_\varphi$ , where  $\Sigma_\Theta$  is defined as follows:

- $\Sigma_\Theta$  includes  $t_i[Z] \rightarrow t_i$  if  $\Theta$  includes a key  $R_i[Z] \rightarrow R_i$ ;
- $\Sigma_\Theta$  includes  $t_i[Z] \subseteq t_j[Z'], t_j[Z'] \rightarrow t_j$  if  $\Theta$  has a foreign key  $R_i[Z] \subseteq R_j[Z'], R_j[Z'] \rightarrow R_j$ .

The set  $\Sigma_\varphi$  consists of the following:

$$D_Y[Y] \rightarrow D_Y, \quad E_X[X] \rightarrow E_X, \quad D_Y[X] \subseteq E_X[X], \quad D_Y[X, Y] \subseteq t_\varphi[X, Y], \quad t_\varphi[XY] \rightarrow t_\varphi,$$

where  $[X, Y]$  stands for the concatenation of list  $X$  and list  $Y$ , and  $t_\varphi$  is the grammar symbol in  $P(R) = t_\varphi^*$ . Observe that  $\text{Att}(R) = X \cup Y$  and thus  $XY$  is a key of  $t_\varphi$ .



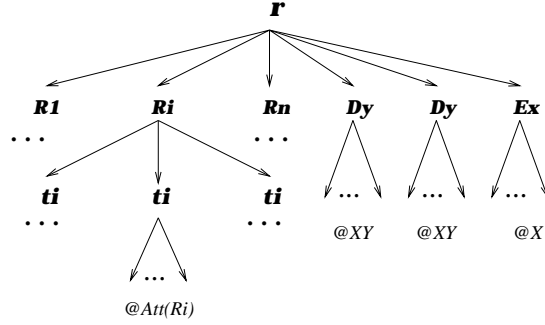


Fig. 2. A tree used in the proof of Theorem 3.1

As depicted in Figure 2, in any XML tree valid w.r.t.  $D$ , there are two distinct  $D_Y$  nodes  $d_1$  and  $d_2$  that have all the attributes in  $X \cup Y$ , and a single  $E_X$  node having all attributes in  $X$ . If  $T \models \Sigma_\varphi$ , then (1)  $d_1[X] = d_2[X]$  by  $D_Y[X] \subseteq E_X[X]$  and the fact  $|ext(E_X)| = 1$ ; and (2)  $d_1[Y] \neq d_2[Y]$  by  $D_Y[Y] \rightarrow D_Y$ . These nodes will serve as a witness for  $\neg\varphi$ .

Given these, we show that  $\bigwedge \Theta \wedge \neg\varphi$  can be satisfied by an instance of  $\mathbf{R}$  if and only if  $\Sigma$  can be satisfied by an XML tree valid w.r.t.  $D$ . Assume that there is an instance  $\mathbf{I}$  of  $\mathbf{R}$  satisfying  $\bigwedge \Theta \wedge \neg\varphi$ . We construct an XML tree  $T$  from  $\mathbf{I}$  as follows. Let  $T$  have a root node  $r$  and a  $R_i$  node for each  $R_i$  in  $\mathbf{R}$ . For any  $R_i \in \mathbf{R}$  and each tuple  $p$  in the instance of  $R_i$  in  $\mathbf{I}$ , we create a distinct  $t_i$  node  $x$  such that  $p.l = x.l$  for all  $l \in Att(R_i)$ . Since  $\mathbf{I} \models \neg\varphi$ , there are two tuples  $p$  and  $p'$  in the instance of  $R$  in  $\mathbf{I}$  such that  $p[X] = p'[X]$  and  $p[Y] \neq p'[Y]$ . We create two distinct  $D_Y$  nodes  $d_1$  and  $d_2$  such that  $d_1.l = p.l$  and  $d_2.l = p'.l$  for all  $l \in Att(R)$ . In addition, we create a single  $E_X$  node  $e$  such that  $e.l = p.l$  for all  $l \in X$ . We define the edge relation of  $T$  such that  $T$  has the form shown in Figure 2. It is easy to verify that  $T \models D$ . By  $\mathbf{I} \models \Theta$  it is easy to verify that  $T \models \Sigma_\Theta$ . By the definition of  $T$ , it is also easy to see that  $T \models \Sigma_\varphi$ . In particular, since  $Att(R) = X \cup Y$  and the set of all attributes of a relation is a key of the relation, we have  $T \models t_\varphi[XY] \rightarrow t_\varphi$ , where  $t_\varphi$  is the symbol in  $P(R) = t_\varphi^*$ . Therefore,  $T \models \Sigma$ . Conversely, suppose that  $D$  has a valid XML tree  $T$  that satisfies  $\Sigma$ . We define an instance  $\mathbf{I}$  of schema  $\mathbf{R}$  as follows. For each  $t_i$  node  $x$ , let  $(l_1 = x.l_1, \dots, l_m = x.l_m)$  be a tuple in the instance of  $R_i$  in  $\mathbf{I}$ , where  $l_1, \dots, l_m$  are an enumeration of  $Att(R_i)$ . Obviously  $\mathbf{I}$  is an instance of  $\mathbf{R}$ . By  $T \models \Sigma_\Theta$ , it is easy to verify that  $\mathbf{I} \models \Theta$ . Moreover, by  $T \models \Sigma_\varphi$  and the definition of  $\mathbf{I}$ , we have  $\mathbf{I} \models \neg\varphi$  since there must be two tuples  $d_1$  and  $d_2$  in the instance of  $R$  in  $\mathbf{I}$  such that  $d_1[X] = d_2[X]$  but  $d_1[Y] \neq d_2[Y]$ . Thus the encoding is indeed a reduction from the complement of the implication problem for keys by keys and foreign keys.

This completes the proof of Theorem 3.1. □

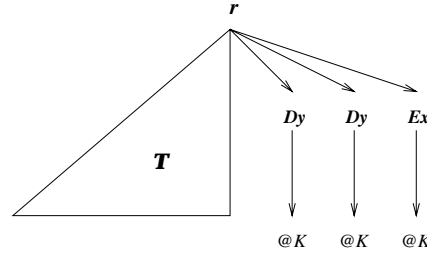


Fig. 3. A tree used in the proof of Lemma 3.3

### 3.2 Undecidability of implication

We next consider the implication problem.

**LEMMA 3.3.** *The following problems are undecidable: given any DTD  $D$ , any set  $\Sigma$  of  $\mathcal{C}_{K,FK}$  constraints over  $D$ , any unary key  $\varphi_1$  and unary inclusion constraint  $\varphi_2$  over  $D$ , whether (1)  $(D, \Sigma) \vdash \varphi_1$ ; (2)  $(D, \Sigma) \vdash \varphi_2$ .*

*Proof:* It suffices to establish a reduction from the consistency problem for  $\mathcal{C}_{K,FK}$  to the complement of the implication problem for  $\mathcal{C}_{K,FK}$ . Let the DTD  $D$  be  $(E, A, P, R, r)$ . We define another DTD  $D' = (E', A', P', R', r)$ , where

$$\begin{aligned}
 E' &= E \cup \{D_Y, E_X\} && \text{where } D_Y, E_X \text{ are fresh element types} \\
 A' &= A \cup \{K\} && \text{where } K \text{ is a fresh attribute} \\
 P'(r) &= P(r), D_Y, D_Y, E_X && \text{i.e., } P(r) \text{ followed by two } D_Y \text{ elements and an } E_X \text{ element} \\
 P'(\tau) &= P(\tau) && \text{for all } \tau \in E \setminus \{r\} \\
 P'(D_Y) &= P'(E_X) = \epsilon \\
 R'(D_Y) &= \{K\} \\
 R'(E_X) &= \{K\} \\
 R'(\tau) &= R(\tau) && \text{for all } \tau \in E
 \end{aligned}$$

We define a unary key  $\varphi_1$ , a unary inclusion constraint  $\varphi_2$  and another key  $\phi$  over  $D'$  as follows:

$$\varphi_1 = D_Y.K \rightarrow D_Y, \quad \varphi_2 = D_Y.K \subseteq E_X.K, \quad \phi = E_X.K \rightarrow E_X.$$

Clearly,  $\Sigma$  is also a set of  $\mathcal{C}_{K,FK}$  constraints over  $D'$ . We next show that (1)  $\Sigma$  is satisfiable over  $D$  iff  $\bigwedge \Sigma \wedge \phi \wedge \varphi_2 \wedge \neg \varphi_1$  is satisfiable over  $D'$ ; (2)  $\Sigma$  is satisfiable over  $D$  iff  $\bigwedge \Sigma \wedge \phi \wedge \varphi_1 \wedge \neg \varphi_2$  is satisfiable over  $D'$ . For if these hold, then the encoding is a reduction from the consistency problem for  $\mathcal{C}_{K,FK}$  to the complements of the implication problems described in Lemma 3.3.

We prove (1) as follows. If there exists a tree  $T \models D'$  and  $T \models \bigwedge \Sigma \wedge \phi \wedge \varphi_2 \wedge \neg \varphi_1$ , then we construct another tree  $T'$  by removing  $D_Y, E_X$  elements from  $T$ . Obviously,  $T' \models D$  and

$T' \models \Sigma$ . Conversely, suppose that there is a tree  $T \models D$  and  $T \models \Sigma$ . We construct another tree  $T'$  from  $T$  as shown in Figure 3. Let us refer to the two  $D_Y$  elements in  $T'$  as  $d_1, d_2$ , and the  $E_X$  element as  $e$ . Let  $d_1.K = d_2.K = e.K$ . Then it is easy to see that  $T' \models D'$ ,  $T' \models \Sigma$  and  $T' \models \phi \wedge \varphi_2 \wedge \neg\varphi_1$ .

We now prove (2). As above, we can show that if there is a tree  $T \models D'$  and  $T \models \bigwedge \Sigma \wedge \phi \wedge \varphi_1 \wedge \neg\varphi_2$ , then there exists another tree  $T'$  such that  $T' \models D$  and  $T' \models \Sigma$ . Conversely, suppose that there is a tree  $T \models D$  and  $T \models \Sigma$ . We construct a tree  $T'$  from  $T$  as shown in Figure 3. Again we refer to the two  $D_Y$  elements in  $T'$  as  $d_1, d_2$ , and the  $E_X$  element as  $e$ . Now let  $d_1.K \neq d_2.K$ . Then it is easy to see that  $T' \models D'$ ,  $T' \models \Sigma$  and  $T' \models \phi \wedge \varphi_1 \wedge \neg\varphi_2$ .  $\square$

From Lemma 3.3 we immediately obtain:

**COROLLARY 3.4.** *For  $\mathcal{C}_{K,FK}$  constraints, the implication problem is undecidable.*

### 3.3 PTIME decidable cases

While the general consistency and implication problems are undecidable, it is possible to identify some decidable cases of low complexity. The first one is checking whether a DTD has a valid XML tree. This is a special case of the consistency problem, namely, when the given set of  $\mathcal{C}_{K,FK}$  constraints is empty. A more interesting special case involves keys only. Let  $\mathcal{C}_K$  denote the set of all keys in  $\mathcal{C}_{K,FK}$ . The *consistency problem for  $\mathcal{C}_K$*  is to determine, given any DTD  $D$  and any set  $\Sigma$  of keys in  $\mathcal{C}_K$  over  $D$ , whether there exists an XML tree valid w.r.t.  $D$  and satisfying  $\Sigma$ . Similarly, we consider the *implication problem for  $\mathcal{C}_{K,FK}$* : given any DTD  $D$ , any set  $\Sigma$  and  $\varphi$  of keys in  $\mathcal{C}_K$  over  $D$ , whether  $(D, \Sigma) \vdash \varphi$ . The next theorem tells that all these cases are decidable.

**THEOREM 3.5.** *The following problems are decidable in linear time:*

- (1) *Given any DTD  $D$ , whether there exists an XML tree valid w.r.t.  $D$ .*
- (2) *The consistency problem for  $\mathcal{C}_K$ .*
- (3) *The implication problem for  $\mathcal{C}_K$ .*

*Proof:* (1) The first problem of the theorem can be reduced to the emptiness problem for a context free grammar (CFG). Observe that a DTD  $D = (E, A, P, R, r)$  can be viewed as an extended CFG  $G_D$  with  $r$  as its start symbol,  $S$  as a nonterminal with a production rule, say,  $S \rightarrow 0$ , and with attributes ( $A$  and  $R$ ) ignored. It is easy to verify that  $D$  has a valid XML tree if and only if  $G_D$  is nonempty, i.e., it generates a terminal string (equivalently, a parse tree). Indeed, given an XML tree  $T$  valid w.r.t.  $D$ , one can construct a parse tree of  $G_D$  by

modifying  $T$ , i.e., by removing attributes from  $T$  and modifying its text nodes. Conversely, given a parse tree  $T'$  of  $G_D$  one can construct a valid XML tree of  $D$  by modifying  $T'$ , i.e., by adding attributes to  $T'$  and removing children of  $\mathbb{S}$  nodes from  $T'$ . It is straightforward to convert the extended CFG  $G_D$  to a CFG  $G$  in linear time, by introducing new nonterminals and their (recursive) production rules to represent Kleene closures. Moreover,  $G_D$  is nonempty if and only if  $G$  is nonempty. It is well known that the emptiness problem for a CFG can be determined in linear time (cf. [24]). Putting everything together, a linear algorithm for checking the validity of  $D$  works as follows: it first generates in linear time the CFG  $G$  from  $D$ , and then checks in linear time whether  $G$  is empty; it concludes that  $D$  has a valid XML tree if and only if  $G$  is nonempty. Thus the validity of DTDs can be decided in linear time.

(2) We next prove the second statement of Theorem 3.5. That is, the consistency problem for  $\mathcal{C}_K$  is decidable in linear time. Given any DTD  $D$  and any set  $\Sigma$  of keys in  $\mathcal{C}_K$  over  $D$ , it suffices to show that  $\Sigma$  can be satisfied by an XML tree valid w.r.t.  $D$  if and only if  $D$  has a valid XML tree. For if it holds, then the second statement follows immediately from the first statement of Theorem 3.5.

We now show the claim. Suppose that there exists an XML tree  $T_1 = (V, lab, ele, att, val, root)$  valid w.r.t.  $D$ . We construct another XML tree  $T_2$  by modifying the  $val$  function in  $T_1$  such that for any key  $\tau[X] \rightarrow \tau$  in  $\Sigma$ ,  $|ext(\tau)| = |ext(\tau.l)|$  in  $T_2$  for every  $l \in X$ . That is,  $T_2 \models \tau.l \rightarrow \tau$  for all  $l \in X$ . More specifically, let  $T_2 = (V, lab, ele, att, val', root)$ . Observe that the only difference between  $T_1$  and  $T_2$  is the definition of the function  $val'$ . For any  $v_1, v_2$  in  $V$  with  $lab(v_1) = \tau$  and  $lab(v_2) = \tau$ , we can make  $val'(att(v_1, l)) \neq val'(att(v_2, l))$  for any  $l \in X$ . Let  $val'(v) = val(v)$  for all other vertices in  $V$ . It is easy to verify that  $T_2$  is valid w.r.t.  $D$  since  $T_1$  is valid w.r.t.  $D$ . In addition,  $T_2 \models \tau[X] \rightarrow \tau$  since for any  $x, y \in ext(\tau)$ ,  $x[X] \neq y[X]$ . The other direction is immediate.

(3) Finally, we prove the last statement of Theorem 3.5. That is, the implication problem for  $\mathcal{C}_K$  is decidable in linear time. To show this, we need the following lemma.

**LEMMA 3.6.** *For any DTD  $D$  and element type  $\tau$  in  $D$ , it is decidable in linear time whether there is an XML tree  $T$  such that  $T \models D$  and moreover,  $|ext(\tau)| > 1$  in  $T$ .*

*Proof:* As in the proof of the first statement of the theorem, it is easy to show that given a DTD  $D$ , one can find in linear time a CFG  $G$  such that  $D$  has a valid XML tree in which  $|ext(\tau)| > 1$  if and only if the start symbol  $r$  of  $G$  derives a terminal string  $w$  whose parse tree has at least two  $\tau$  nodes. This can be transformed in linear time to the problem of checking if a given CFG derives a string with at least two occurrences of a given terminal symbol, which in turn can be solved in linear time by a minor modification of the emptiness test for CFG from [24].  $\square$

Let  $\Sigma$  be a set of keys in  $\mathcal{C}_K$  over  $D$ , and  $\varphi = \tau[X] \rightarrow \tau$  be another key in  $\mathcal{C}_K$  over  $D$ . We say that  $\Sigma$  *subsumes*  $\varphi$  if there is  $\phi = \tau[Y] \rightarrow \tau$  in  $\Sigma$  such that  $Y \subseteq X$ , i.e.,  $\varphi$  is a *superkey* of  $\phi$ . Using this and Lemma 3.6 we can prove the following:

**LEMMA 3.7.** *Let  $D$  be a DTD,  $\Sigma$  a set of keys in  $\mathcal{C}_K$  over  $D$ , and  $\varphi = \tau[X] \rightarrow \tau$  another key in  $\mathcal{C}_K$  over  $D$ . There is an XML tree  $T$  such that  $T \models D$ ,  $T \models \Sigma$  but  $T \models \neg\varphi$  if and only if  $\Sigma$  does not subsume  $\varphi$  and moreover, there is an XML tree  $T'$  such that  $T' \models D$  and  $|ext(\tau)| > 1$  in  $T'$ . In addition, this is decidable in linear time in the sizes of  $D$  and  $\Sigma \cup \{\varphi\}$ .*

*Proof:* We first show that there is an XML tree  $T$  such that  $T \models D$ ,  $T \models \Sigma$  but  $T \models \neg\varphi$  iff  $\Sigma$  does not subsume  $\varphi$  and moreover, there is an XML tree  $T'$  such that  $T' \models D$  and  $|ext(\tau)| > 1$  in  $T'$ . Suppose that there is an XML tree  $T$  such that  $T \models D$ ,  $T \models \Sigma$  and  $T \models \neg\varphi$ . Then obviously,  $T$  is valid w.r.t.  $D$ , and moreover, there must be at least two  $\tau$  elements  $d_1, d_2$  in  $T$  such that  $d_1[X] = d_2[X]$  but  $d_1 \neq d_2$  since  $T \models \neg\varphi$ . Thus there must be  $|ext(\tau)| > 1$  in  $T$ . In addition,  $\Sigma$  cannot contain  $\tau[Y] \rightarrow \tau$  with  $Y \subseteq X$ , since otherwise it would contradict  $T \models \neg\varphi$  and  $T \models \Sigma$ . Conversely, let  $T'$  be a tree such that  $T' \models D$  and  $|ext(\tau)| > 1$  in  $T'$ . Thus there are at least two  $\tau$  elements  $d_1, d_2$  in  $T'$ . We construct a new tree  $T$  by modifying the string values associated with the attributes of  $T'$ , while leaving the other functions of  $T'$  unchanged. More specifically, we let  $d_1[X] = d_2[X]$  in  $T$  but all other attributes have different string values. It is easy to verify that  $T \models D$  and  $T \models \neg\varphi$  by the definition of  $T$ . To show  $T \models \Sigma$ , suppose by contradiction that there were  $\phi \in \Sigma$  such that  $T \models \neg\phi$ . Then  $\phi$  must be of the form  $\tau[Y] \rightarrow \tau$  where  $Y \subseteq X$ , i.e.,  $\varphi$  is a superkey of  $\phi$ , since except  $d_1[X] = d_2[X]$ , distinct nodes in  $T$  have the different attribute values by the definition of  $T$ . This contradicts the assumption that  $\Sigma$  does not subsume  $\varphi$ . Thus the first statement of the lemma holds.  $\square$

To show that this can be done in linear time, observe that by Lemma 3.6, it can be decided in linear time in the size of  $D$  whether there is a tree  $T$  such that  $T \models D$  and  $|ext(\tau)| > 1$  in  $T$ . In addition, it is decidable in linear time in the size of  $\Sigma \cup \{\varphi\}$  whether  $\varphi$  is a superkey of some key in  $\Sigma$  (see e.g., [1] for discussions about a linear time algorithm for checking implication of functional dependencies). Thus it is decidable in linear time in the sizes of  $D$  and  $\Sigma \cup \{\varphi\}$  whether these conditions hold.  $\square$

This suffices to prove the third statement of Theorem 3.5 because  $(D, \Sigma) \vdash \varphi$  iff there is no XML tree  $T$  such that  $T \models D$ ,  $T \models \Sigma$  but  $T \models \neg\varphi$ . By Lemma 3.7, the latter can be decided in linear time.

This completes the proof of Theorem 3.5.  $\square$

Given Theorem 3.5, one would be tempted to think that when only foreign keys are considered, the analyses of consistency and implication could also be simpler. However, it is not the case.

Recall that a foreign key of  $\mathcal{C}_{K,FK}$  consists of an inclusion constraint and a key. Thus we cannot exclude keys in the presence of foreign keys. It is not hard to show that consistency and implication of foreign keys in  $\mathcal{C}_{K,FK}$  remain undecidable.

#### 4. UNARY KEYS AND FOREIGN KEYS

The undecidability of the consistency problem for general keys and foreign keys motivates us to look for restricted classes of constraints. One important class is  $\mathcal{C}_{K,FK}^{Unary}$ , the class of unary keys and foreign keys. A cursory examination of existing XML specifications reveals that most keys and foreign keys are single-attribute constraints, i.e., unary. In particular, in XML DTDs, one can only specify unary constraints with ID and IDREF attributes.

In this section, we first investigate the consistency problem for  $\mathcal{C}_{K,FK}^{Unary}$ . To simplify the discussion and to establish a (slightly) stronger result, we consider a larger class of constraints, namely,  $\mathcal{C}_{K,IC}^{Unary}$ , the class of unary keys and unary inclusion constraints. In contrast to  $\mathcal{C}_{K,FK}^{Unary}$ ,  $\mathcal{C}_{K,IC}^{Unary}$  allows the presence of unary inclusion constraints independent of keys. We develop an encoding of DTDs and  $\mathcal{C}_{K,IC}^{Unary}$  constraints with linear integer constraints. This enables us to reduce the consistency problem for  $\mathcal{C}_{K,IC}^{Unary}$  (and thus for  $\mathcal{C}_{K,FK}^{Unary}$ ) to the linear integer programming problem, one of the most studied NP-complete problems. We then use the same technique to show that the consistency problem remains in NP when negations of keys are allowed, i.e., the problem for  $\mathcal{C}_{K,IC}^{Unary}$  constraints is also in NP. Finally, we identify several tractable cases of the consistency problems.

##### 4.1 Coding DTDs, unary constraints

We show that  $\mathcal{C}_{K,IC}^{Unary}$  constraints and DTDs can be encoded with linear equalities and inequalities on the integers, called *cardinality constraints*. The encoding allows us to reduce the consistency problem for  $\mathcal{C}_{K,IC}^{Unary}$  constraints in PTIME to the *linear integer programming (LIP)* problem:

##### LINEAR INTEGER PROGRAMMING (LIP)

INPUT: An  $m \times n$  matrix  $A$  of integers and a column vector  $\vec{b}$  of  $m$  integers.

QUESTION: Does there exist a column vector  $\vec{x}$  of  $n$  integers such that  $A\vec{x} \geq \vec{b}$ ?

That is, for  $i \in [1, m]$ ,

$$\sum_{j \in [1, n]} a_{ij} x_j \geq b_i,$$

where  $a_{ij}$  is the  $j$ th element of the  $i$ th row of  $A$ ,  $x_j$  is the  $j$ th entry of  $\vec{x}$  and  $b_i$  is the  $i$ th entry of  $\vec{b}$ . It is known that LIP is NP-complete in the strong sense [22]. In particular, when nonnegative integer solutions are considered, [32] has shown that if the problem has a solution, then it has another solution in which for all  $j \in [1, n]$ ,  $x_j$  is no larger than  $n(ma)^{2m+1}$ , where  $a$  is the largest absolute value of elements in  $A$  and  $\vec{b}$ .

More specifically, we show the following:

**THEOREM 4.1.** *There is a polynomial ( $O(s^2 \cdot \log s)$ ) time algorithm that, given a DTD  $D$  and a set  $\Sigma$  of  $\mathcal{C}_{K,IC}^{Unary}$  constraints, constructs an integer matrix  $A$  and an integer vector  $\vec{b}$  such that there exists an XML tree valid w.r.t.  $D$  and satisfying  $\Sigma$  if and only if  $A\vec{x} \geq \vec{b}$  has an integer solution.*

As an immediate result, we have:

**COROLLARY 4.2.** *The consistency problem for  $\mathcal{C}_{K,FK}^{Unary}$  constraints is in NP.*

The proof of Theorem 4.1 is a bit involved. A road map of the proof is as follows. Given a DTD  $D$  and a set  $\Sigma$  of  $\mathcal{C}_{K,IC}^{Unary}$  constraints over  $D$ , we define in  $O(s^2 \cdot \log s)$  time (in the sizes of  $D$  and  $\Sigma$ , denoted by  $|D|$  and  $|\Sigma|$ , respectively) the following:

- another DTD  $D_N$ , referred to as the *simplified DTD of  $D$* , in which regular expressions are restricted to have at most one operator: either “|” (union) or “,” (concatenation)<sup>2</sup>; we reduce the consistency of  $D$  and  $\Sigma$  to that of  $D_N$  and  $\Sigma$ , i.e., there exists an XML tree valid w.r.t.  $D$  and satisfying  $\Sigma$  if and only if there exists an XML tree valid w.r.t.  $D_N$  and satisfying  $\Sigma$ ;
- a set  $C_\Sigma$  of linear integer constraints such that there is an XML tree valid w.r.t.  $D_N$  and satisfying  $\Sigma$  if and only if there is an XML tree valid w.r.t.  $D_N$  and satisfying  $C_\Sigma$ ;
- a system  $\Psi_{D_N}$  of linear integer constraints such that there exists an XML tree valid w.r.t.  $D_N$  if and only if  $\Psi_{D_N}$  admits an integer solution; the cardinality constraints in  $\Psi_{D_N}$  are more complex than those studied in the context of object-oriented and relational databases [12; 13; 26];
- finally, a system of integer constraints  $\Psi(D, \Sigma)$  from  $C_\Sigma$  and  $\Psi_{D_N}$  such that there exists an XML tree valid w.r.t.  $D$  and satisfying  $\Sigma$  if and only if  $\Psi(D, \Sigma)$  admits an integer solution.

Putting everything together, we reduce the consistency problem for  $\mathcal{C}_{K,IC}^{Unary}$  to the existence of a solution of an instance of LIP, and thus obtain the NP bound.

<sup>2</sup>We are grateful to one of the referees for suggesting this simplification of DTDs.

*Proof of Theorem 4.1:* We start by describing the process of simplifying DTDs. We shall then present an encoding of unary constraints and DTDs. Finally, we develop a characterization of XML specifications with both DTDs and unary constraints in terms of linear integer constraints.

**Simplifying DTDs.** We first explain how to reduce the consistency problem for  $\mathcal{C}_{K,IC}^{Unary}$  to that over simple DTDs. Intuitively, we replace long regular expressions in  $P(\tau)$  by shorter ones. Formally, consider a DTD  $D = (E, A, P, R, r)$ . For each  $\tau \in E$ ,  $P(\tau)$  is a regular expression  $\alpha$ . A DTD is basically an extended regular grammar (cf. [11; 31]); thus  $\tau \rightarrow \alpha$  can be viewed as the production rule for  $\tau$ . We rewrite the regular expression  $\alpha$  by introducing a set  $N$  of *new* element types (nonterminals) such that the production rules of the new DTD have one of the following forms:

$$\tau \rightarrow \tau_1, \tau_2 \quad \tau \rightarrow \tau_1 \mid \tau_2 \quad \tau \rightarrow \tau_1 \quad \tau \rightarrow \mathbf{S} \quad \tau \rightarrow \epsilon$$

where  $\tau, \tau_1, \tau_2$  are element types in  $E \cup N$ ,  $\mathbf{S}$  is the string type and  $\epsilon$  denotes the empty word. More specifically, we conduct the following “simplifying” process on the production rule  $\tau \rightarrow \alpha$ :

- (1) If  $\alpha = (\alpha_1, \alpha_2)$ , then we introduce two new element types  $\tau_1, \tau_2$  and replace  $\tau \rightarrow \alpha$  with a new rule  $\tau \rightarrow \tau_1, \tau_2$ . We proceed to process  $\tau_1 \rightarrow \alpha_1$  and  $\tau_2 \rightarrow \alpha_2$  in the same way.
- (2) If  $\alpha = (\alpha_1 \mid \alpha_2)$ , then we introduce two new element types  $\tau_1, \tau_2$  and replace  $\tau \rightarrow \alpha$  with a new rule  $\tau \rightarrow \tau_1 \mid \tau_2$ . We proceed to process  $\tau_1 \rightarrow \alpha_1$  and  $\tau_2 \rightarrow \alpha_2$  in the same way.
- (3) If  $\alpha = \alpha_1^*$ , then we introduce a new element type  $\tau_1$  and replace  $\tau \rightarrow \alpha$  with  $\tau \rightarrow \tau_1$ . We proceed to process  $\tau_1 \rightarrow \epsilon \mid \alpha_1, \tau_1$  in the same way.
- (4) If  $\alpha$  is one of  $\tau' \in E$ ,  $\mathbf{S}$  or  $\epsilon$ , then the rule for  $\tau$  remains unchanged.

To avoid introducing unnecessary new element types, in the first two cases above, if  $\alpha_1$  (resp.  $\alpha_2$ ) is a symbol of  $E \cup \{\mathbf{S}\}$ , we do not introduce a new element type for  $\alpha_1$  (resp.  $\alpha_2$ ).

We refer to the set of new element types introduced when processing  $\tau \rightarrow P(\tau)$  as  $N_\tau$  and the set of production rules generated/revised as  $P_\tau$ . Note that  $N_\tau \cap E = \emptyset$  for any  $\tau \in E$ .

We define a new DTD  $D_N = (E_N, A, P_N, R_N, r)$ , referred to as the *simplified DTD* of  $D$  (or just a *simple* DTD if  $D$  is clear from the context), where

- $E_N = E \cup \bigcup_{\tau \in E} N_\tau$ , i.e.,  $E$  plus those new element types introduced in the simplifying process;
- $P_N = \bigcup_{\tau \in E} P_\tau$ , i.e., production rules generated/revised in the simplifying process;
- $R_N(\tau) = R(\tau)$  for each  $\tau \in E$ , and  $R_N(\tau) = \emptyset$  for each  $\tau \in E_N \setminus E$ .



Note that the root element type  $r$  and the set  $A$  of attributes remain unchanged. Moreover, elements of any type in  $E_N \setminus E$  do not have any attribute. Note that  $D_N$  does not contain the Kleene star “\*”.

For example, the simplified DTD of  $D_1$  given in Section 1 is  $D_1^N = (E_1^N, A_1, P_1^N, R_1^N, r)$ , where

$$\begin{aligned}
E_1^N &= \{teachers, teacher, teach, research, subject, \tau_t^1, \tau_t^2, \tau_\epsilon\} \\
A_1 &= \{name, taught\_by\} \\
P_1^N(teachers) &= teacher, \tau_t^1 \\
P_1^N(\tau_t^1) &= \tau_\epsilon \mid \tau_t^2 \\
P_1^N(\tau_\epsilon) &= \epsilon \\
P_1^N(\tau_t^2) &= teacher, \tau_t^1 \\
P_1^N(teacher) &= teach, research \\
P_1^N(teach) &= subject, subject \\
P_1^N(subject) &= P_1^N(research) = \mathbb{S} \\
R_1^N(teacher) &= \{name\} \\
R_1^N(subject) &= \{taught\_by\} \\
R_1^N(teachers) &= R_1^N(teach) = R_1^N(research) = R_1^N(\tau_t^1) = R_1^N(\tau_t^2) = R_1^N(\tau_\epsilon) = \emptyset \\
r_1 &= teachers
\end{aligned}$$

Here  $\tau_t^1, \tau_t^2, \tau_\epsilon$  are the new element types introduced.

The simplified DTD  $D_2^N$  of  $D_2$  in Section 1 is the same as  $D_2$  itself.

Obviously, any set  $\Sigma$  of  $\mathcal{C}_{K,IC}^{Unary}$  constraints over  $D$  is also a set of  $\mathcal{C}_{K,IC}^{Unary}$  constraints over the simplified DTD  $D_N$  of  $D$ . The next lemma establishes the connection between  $D$  and  $D_N$ , which allows us to consider only simple DTDs from now on.

**LEMMA 4.3.** *Let  $D$  be a DTD,  $D_N$  be the simplified DTD of  $D$  and  $\Sigma$  be a set of  $\mathcal{C}_{K,IC}^{Unary}$  constraints over  $D$ . Then there exists an XML tree  $T_1$  such that  $T_1 \models D$  and  $T_1 \models \Sigma$  iff there exists an XML tree  $T_2$  such that  $T_2 \models D_N$  and  $T_2 \models \Sigma$ .*

*Proof:* It suffices to show the following claim. For any XML tree  $T_1 \models D$  one can construct an XML tree  $T_2 \models D_N$ , and for any  $T_2 \models D_N$  one can construct  $T_1 \models D$ , such that for any element type  $\tau$  in  $D$  and  $l \in R(\tau)$ ,  $|ext(\tau)|$  in  $T_2$  equals  $|ext(\tau)|$  in  $T_1$ , and  $ext(\tau.l)$  in  $T_2$  equals  $ext(\tau.l)$  in  $T_1$ .

We first prove the lemma assuming that the claim is true. Assume that there exists an XML tree  $T_1$  such that  $T_1 \models D$  and  $T_1 \models \Sigma$ . Find the tree  $T_2 \models D_N$  as in the claim. Suppose that

there is  $\varphi \in \Sigma$  such that  $T_2 \not\models \varphi$ . If  $\varphi$  is a key  $\tau.l \rightarrow \beta.\tau$ , then there are two distinct nodes  $x, y \in \text{ext}(\tau)$  in  $T_1$  such that  $x.l = y.l$ . Thus  $|\text{ext}(\tau.l)| < |\text{ext}(\tau)|$  in  $T_2$  since every  $\tau$  element has a single  $l$  attribute. Since  $T_1 \models \varphi$ , it must be the case that  $|\text{ext}(\tau.l)| = |\text{ext}(\tau)|$  in  $T_1$  since the value  $x.l$  of each  $x \in \text{ext}(\tau)$  uniquely identifies  $x$  among all the nodes in  $\text{ext}(\tau)$ . This contradicts the claim that  $|\text{ext}(\tau)|$  in  $T_2$  equals  $|\text{ext}(\tau)|$  in  $T_1$  and  $\text{ext}(\tau.l)$  in  $T_2$  equals  $\text{ext}(\tau.l)$  in  $T_1$ . If  $\varphi$  is an inclusion constraint  $\tau_1.l_1 \subseteq \tau_2.l_2$ , then there is  $x \in \text{ext}(\tau_1)$  such that for all  $y \in \text{ext}(\tau_2)$  in  $T_2$ ,  $x.l_1 \neq y.l_2$ . That is,  $x.l_1 \notin \text{ext}(\tau_2.l_2)$ . By the claim,  $x.l_1 \in \text{ext}(\tau_1.l_1)$  in  $T_1$ . Since  $T_1 \models \varphi$ , we have  $x.l_1 \in \text{ext}(\tau_2.l_2)$  in  $T_1$ . Again by the claim, we have  $x.l_1 \in \text{ext}(\tau_2.l_2)$  in  $T_2$ , which contradicts the assumption. The proof for the other direction is similar.

We next verify the claim. Given an XML tree  $T_1 = (V_1, \text{lab}_1, \text{ele}_1, \text{att}, \text{val}, \text{root})$  such that  $T_1 \models D$ , we construct an XML tree  $T_2$  by modifying  $T_1$  such that  $T_2 \models D_N$ . Consider a  $\tau$  element  $v$  in  $T_1$ . Let  $\text{ele}_1(v) = [v_1, \dots, v_n]$  and  $w = \text{lab}_1(v_1) \dots \text{lab}_1(v_n)$ . Recall  $N_\tau$  and  $P_\tau$ , the set of nonterminals and the set of production rules generated when simplifying  $\tau \rightarrow P(\tau)$ . Let  $Q_\tau$  be the set of  $E$  symbols that appear in  $P_\tau$  plus  $S$ . We can view  $G = (Q_\tau, N_\tau \cup \{\tau\}, P_\tau, \tau)$  as a context free grammar, where  $Q_\tau$  is the set of terminals,  $N_\tau \cup \{\tau\}$  the set of nonterminals,  $P_\tau$  the set of production rules and  $\tau$  the start symbol. Since  $T_1 \models D$ , we have  $w \in P(\tau)$ . By a straightforward induction on the structure of  $P_N(\tau)$  it can be verified that  $w$  is in the language defined by  $G$ . Thus there is a parse tree  $T(w)$  of the grammar  $G$  for  $w$ , and  $w$  is the frontier (the list of leaves from left to right) of  $T(w)$ . Without loss of generality, assume that the root of  $T(w)$  is  $v$ , and the leaves are  $v_1, \dots, v_n$ . Intuitively, we construct  $T_2$  by replacing each element  $v$  in  $T_1$  by such a parse tree. More specifically, let  $T_2 = (V_2, \text{lab}_2, \text{ele}_2, \text{att}, \text{val}, \text{root})$ . Here  $V_2$  consists of nodes in  $V_1$  and the internal nodes introduced in the parse trees. For each  $x$  in  $V_2$ , let  $\text{lab}_2(x) = \text{lab}_1(x)$  if  $x \in V_1$ , and otherwise let  $\text{lab}_2(x)$  be the node label of  $x$  in the parse tree where  $x$  belongs. Note that nodes in  $V_2 \setminus V_1$  are elements of some type in  $E_N \setminus E$ . If  $\text{lab}_2(x)$  is an element type, let  $\text{ele}_2(x)$  be the list of its children in the parse tree. Note that  $\text{att}$  and  $\text{val}$  remain unchanged. By the construction of  $T_2$  it can be verified that  $T_2 \models D_N$ . Moreover, for any  $\tau \in E$  and  $l \in R(\tau)$ ,  $|\text{ext}(\tau)|$  in  $T_2$  equals  $|\text{ext}(\tau)|$  in  $T_1$  and  $\text{ext}(\tau.l)$  in  $T_2$  equals  $\text{ext}(\tau.l)$  in  $T_1$  because none of the new nodes, i.e., nodes in  $V_2 \setminus V_1$ , is labeled with an  $E$  type, and the function  $\text{att}$  remains unchanged.

Conversely, assume that there is  $T_2 = (V_2, \text{lab}_2, \text{ele}_2, \text{att}, \text{val}, \text{root})$  such that  $T_2 \models D_N$ . We construct  $T_1$  by modifying  $T_2$  such that  $T_1 \models D$ . For any node  $v \in V_2$  with  $\text{lab}(v) = \tau$  and  $\tau \in E_N \setminus E$ , we substitute the subelements of  $v$  for  $v$  in  $\text{ele}(v')$ , where  $v'$  is the parent of  $v$ . In addition, we remove  $v$  from  $V_2$ ,  $\text{lab}_2(v)$  from  $\text{lab}_2$ , and  $\text{ele}_2(v)$  from  $\text{ele}_2$ . Observe that by the definition of  $D_N$ , no attributes are defined for elements of any type in  $E_N \setminus E$ . We repeat the process until there is no node labeled with element type in  $E_N \setminus E$ . Now let  $T_1 = (V_1, \text{lab}_1, \text{ele}_1, \text{att}, \text{val}, \text{root})$ , where  $V_1$ ,  $\text{lab}_1$  and  $\text{ele}_1$  are  $V_2$ ,  $\text{lab}_2$  and  $\text{ele}_2$  at the end of the process, respectively. Observe that  $\text{att}$ ,  $\text{val}$  and  $\text{root}$  remain unchanged. By the definition of  $T_1$  it can be verified that  $T_1 \models D$ ; and in addition, for any  $\tau \in E$  and  $l \in R(\tau)$ ,  $|\text{ext}(\tau)|$

in  $T_1$  equals  $|ext(\tau)|$  in  $T_2$ , and  $ext(\tau.l)$  in  $T_1$  equals  $ext(\tau.l)$  in  $T_2$ , because none of the nodes removed is labeled with a type of  $E$  and the functions  $att$  and  $val$  are unchanged.  $\square$

It is easy to see that  $D_N$  is computable in linear time in the size of  $D$ .

**Encoding unary constraints.** We now give a coding of  $\mathcal{C}_{K,IC}^{Unary}$  constraints. Let  $\Sigma$  be a set of  $\mathcal{C}_{K,IC}^{Unary}$  constraints over DTD  $D$  and  $D_N$  be simplified DTD of  $D$ . Referring to an arbitrary XML tree  $T$  valid w.r.t.  $D$ , we derive from  $\Sigma$  a class of linear integer constraints on  $T$ , denoted by  $C_\Sigma$  and referred to *the cardinality constraints determined by  $\Sigma$* , as follows. For any  $\varphi \in \Sigma$ ,

- if  $\varphi$  is a key constraint  $\tau.l \rightarrow \tau$ , then  $|ext(\tau.l)| = |ext(\tau)|$  is in  $C_\Sigma$ ;
- if  $\varphi$  is an inclusion constraint  $\tau_1.l_1 \subseteq \tau_2.l_2$ , then  $|ext(\tau_1.l_1)| \leq |ext(\tau_2.l_2)|$  is in  $C_\Sigma$ .
- $|ext(\tau.l)| \leq |ext(\tau)|$  and  $0 \leq |ext(\tau.l)|$  are in  $C_\Sigma$  for any  $\tau \in E$  and  $l \in R(\tau)$ .

We use  $T \models C_\Sigma$  to denote that  $T$  satisfies all constraints of  $C_\Sigma$ .

For example, recall the set  $\Sigma_1$  of  $\mathcal{C}_{K,FK}^{Unary}$  constraints over the DTD  $D_1$  given in Section 1. The set of cardinality constraints determined by  $\Sigma_1$ , denoted by  $C_{\Sigma_1}$ , consists of:

$$\begin{aligned} |ext(teacher.name)| &= |ext(teacher)| \\ |ext(subject.taught\_by)| &= |ext(subject)| \\ |ext(subject.taught\_by)| &\leq |ext(teacher.name)| \\ 0 &\leq |ext(teacher.name)| \\ 0 &\leq |ext(subject.taught\_by)| \end{aligned}$$

It is worth mentioning that  $|ext(\tau.l)| = |ext(\tau)|$  characterizes a key  $\tau.l \rightarrow \tau$ . Indeed, for any XML tree  $T$  valid w.r.t.  $D_N$ ,  $T \models |ext(\tau.l)| = |ext(\tau)|$  iff  $T \models \tau.l \rightarrow \tau$ . However, things can go wrong when it comes to inclusion constraints. Although  $T \models \tau_1.l_1 \subseteq \tau_2.l_2$  implies  $T \models |ext(\tau_1.l_1)| \leq |ext(\tau_2.l_2)|$ , the other direction does not necessarily hold. This does not lose generality as we do not intend to capture negations of inclusion constraints with this coding. Indeed, the lemma below shows that we are able to consider  $C_\Sigma$  instead of  $\Sigma$  when studying the consistency of  $\Sigma$ .

**LEMMA 4.4.** *Let  $D_N$  be a simplified DTD of  $D$ ,  $\Sigma$  be a set of  $\mathcal{C}_{K,IC}^{Unary}$  constraints over  $D$ , and  $C_\Sigma$  be the set of cardinality constraints determined by  $\Sigma$ . Then there exists an XML tree  $T_1$  such that  $T_1 \models D_N$  and  $T_1 \models \Sigma$  if and only if there exists an XML tree  $T_2$  such that  $T_2 \models D_N$  and  $T_2 \models C_\Sigma$ . In addition, any XML tree valid w.r.t.  $D_N$  and satisfying  $\Sigma$  also satisfies  $C_\Sigma$ .*

*Proof:* It is easy to see that for any XML tree  $T_1$  that satisfies  $\Sigma$ , it must be the case that  $T_1 \models C_\Sigma$ . Conversely, we show that if there exists an XML tree  $T_2 = (V, lab, ele, att, val, root)$

such that  $T_2 \models D_N$  and  $T_2 \models C_\Sigma$ , then we can construct an XML tree  $T_1$  such that  $T_1 \models D_N$  and  $T_1 \models \Sigma$ .

We construct  $T_1$  from  $T_2$  by modifying the function  $val$  while leaving  $V$ ,  $lab$ ,  $ele$ ,  $att$  and  $root$  unchanged. As cardinality constraints of  $C_\Sigma$  do not involve text nodes, we change  $val$  for attributes only. More specifically, we modify  $val(v)$  if  $lab(v) \in A$ , i.e., if  $v$  is an attribute, and leave  $val(v)$  unchanged otherwise. Let  $S = \{\tau.l \mid \tau \in E, l \in R(\tau)\}$ . To define the new function, denoted by  $val'$ , we first associate a set  $V_{\tau.l}$  of string values with each  $\tau.l$  in  $S$ . Let  $N$  be the maximum cardinality of  $ext(\tau.l)$  in  $T_2$ , i.e.,  $N \geq |ext(\tau.l)|$  in  $T_2$  for all  $\tau.l \in S$ . Let  $V_S = \{a_i \mid i \in [1, N]\}$  be a set of distinct string values. For each  $\tau.l \in S$ , let  $V_{\tau.l} = \{a_i \mid i \in [1, |ext(\tau.l)|]\}$ , and for each  $x \in ext(\tau)$ , let  $val'(att(x, l))$  be a string value in  $V_{\tau.l}$  such that in  $T_1$ ,  $ext(\tau.l) = V_{\tau.l}$ . In addition, for each key  $\tau.l \rightarrow \tau$  in  $\Sigma$ , let  $x.l$  be a distinct string value in  $V_{\tau.l}$ . This is possible because by the definition of  $T_1$ , (1)  $ext(\tau)$  in  $T_1$  equals  $ext(\tau)$  in  $T_2$ ; (2)  $|ext(\tau.l)|$  in  $T_1$  equals  $|ext(\tau.l)|$  in  $T_2$ ; and (3)  $T_2 \models C_\Sigma$  and  $|ext(\tau)| = |ext(\tau.l)|$  is in  $C_\Sigma$ . We next show that  $T_1$  is indeed what we want. It is easy to verify that  $T_1 \models D_N$  given the construction of  $T_1$  from  $T_2$  and the assumption that  $T_2 \models D_N$ . To show that  $T_1 \models \Sigma$ , we consider  $\varphi \in \Sigma$  in the following cases. (1) If  $\varphi$  is a key  $\tau.l \rightarrow \tau$ , it is immediate from the definition of  $T_1$  that  $T_1 \models \varphi$  since for any  $x \in ext(\tau)$ ,  $x.l$  is a distinct string value in  $V_{\tau.l}$ . (2) If  $\varphi$  is  $\tau_1.l_1 \subseteq \tau_2.l_2$ , then  $T_2 \models |ext(\tau_1.l_1)| \leq |ext(\tau_2.l_2)|$  by  $T_2 \models C_\Sigma$ . Recall that by the definition of  $val'$ , for  $i \in [1, 2]$ ,  $V_{\tau_i.l_i} = \{a_i \mid i \in [1, |ext(\tau_i.l_i)|]\}$  and in  $T_1$ ,  $ext(\tau_i.l_i) = V_{\tau_i.l_i}$ . Thus  $ext(\tau_1.l_1) \subseteq ext(\tau_2.l_2)$  in  $T_1$ . That is,  $T_1 \models \varphi$ . Therefore,  $T_1 \models D_N$  and  $T_1 \models \Sigma$ .  $\square$

Observe that in the construction of  $T_1$  above, it is possible that  $ext(\tau_1.l_1) \subseteq ext(\tau_2.l_2)$  even if  $\Sigma$  does not imply  $\tau_1.l_1 \subseteq \tau_2.l_2$ . This does not have an impact on the consistency analysis, as negations of inclusion constraints are not involved in the analysis.

It is straightforward to verify that given any set  $\Sigma$  of  $\mathcal{C}_{K,IC}^{Unary}$  constraints over a DTD  $D$ , the set  $C_\Sigma$  of cardinality constraints determined by  $\Sigma$  can be computed in linear time in  $|\Sigma|$  and  $|D|$ .

**Encoding DTDs.** We next move to a coding of DTDs. By Lemma 4.3 we can consider simple DTDs only. Given any simple DTD  $D = (E, A, P, R, r)$ , we encode it in linear time with a system  $\Psi_D$  of linear integer constraints such that  $D$  has a valid XML tree if and only if  $\Psi_D$  has an integer solution.

We first describe the variables used in the system  $\Psi_D$ . For each symbol  $\tau \in E \cup \{\mathbf{S}\}$ ,  $|ext(\tau)|$  is a distinct variable. Intuitively, in an XML tree  $T$  conforming to  $D$ ,  $|ext(\tau)|$  keeps track of the number of all  $\tau$  elements. In addition, for each occurrence of  $\tau$  in the definition  $P(\tau')$  of some element type  $\tau'$ , we also create a distinct variable. More specifically, we create such variables as follows: if  $P(\tau') = \tau_1$  for  $\tau_1 \in E \cup \{\mathbf{S}\}$ , then we create a distinct variable  $x_{\tau_1, \tau'}^1$ ; if  $P(\tau') = (\tau_1, \tau_2)$  or  $P(\tau') = (\tau_1 | \tau_2)$ , then we create two distinct variables  $x_{\tau_1, \tau'}^1$  and  $x_{\tau_2, \tau'}^2$ .

Intuitively, for  $i \in [1, 2]$ ,  $x_{\tau_1, \tau'}^i$  keeps track of the number of  $\tau_i$  subelements at position  $i$  under all  $\tau'$  elements in  $T$ . For example, given an element type definition  $P(\text{teach}) = \text{subject}, \text{subject}$ , we create two distinct variables  $x_{(\text{subject}, \text{teach})}^1$  and  $x_{(\text{subject}, \text{teach})}^2$ . Let  $X_\tau$  be the set of all variables of the form  $x_{\tau, \tau'}^i$ .

Using these variables, for each  $\tau \in E$ , we define a set  $\psi_\tau$  of linear integer constraints that characterizes  $P(\tau)$  quantitatively, as follows:

- If  $P(\tau) = \tau_1$  for  $\tau_1 \in E \cup \{\mathbf{S}\}$ , then  $\psi_\tau$  includes  $|ext(\tau)| = x_{\tau_1, \tau}^1$ . Referring to the XML tree  $T$ , this assures that each  $\tau$  element has a unique  $\tau_1$  subelement.
- If  $P(\tau) = (\tau_1, \tau_2)$ , then  $\psi_\tau$  includes  $|ext(\tau)| = x_{\tau_1, \tau}^1$  and  $|ext(\tau)| = x_{\tau_2, \tau}^2$ . These assure that each  $\tau$  element in  $T$  must have a unique  $\tau_1$  subelement and a unique  $\tau_2$  subelement.
- If  $P(\tau) = (\tau_1 | \tau_2)$ , then  $\psi_\tau$  includes  $|ext(\tau)| = x_{\tau_1, \tau}^1 + x_{\tau_2, \tau}^2$ . These assure that each  $\tau$  element in  $T$  must have either a  $\tau_1$  subelement or a  $\tau_2$  subelement, and thus the sum of the number of these  $\tau_1$  subelements and the number of  $\tau_2$  subelements equals the number of  $\tau$  elements in  $T$ .

The set of cardinality constraints determined by DTD  $D$ , denoted by  $\Psi_D$ , consists of the following:

- $|ext(r)| = 1$ ; i.e., there is a unique root in any XML tree valid w.r.t.  $D$ ;
- constraints of  $\psi_\tau$  for each  $\tau \in E$ ; these assure that  $P(\tau)$  is satisfied;
- $|ext(\tau)| = \sum_{x_{\tau, \tau'}^i \in X_\tau} x_{\tau, \tau'}^i$  for each  $\tau \in (E \setminus \{r\}) \cup \{\mathbf{S}\}$ ; this indicates that the set  $ext(\tau)$  includes all  $\tau$  elements no matter where they occur in an XML tree;
- $x \geq 0$  for any variable  $x$  used above; i.e., the number of elements (subelements) is nonnegative.

We say that  $\Psi_D$  is *consistent* if and only if  $\Psi_D$  admits an integer solution. That is, there is an integer assignment to the variables of  $\Psi_D$  such that all the linear integer constraints in  $\Psi_D$  are satisfied.

As an example, let us consider the simple DTDs  $D_1^N$  and  $D_2^N$  given above. The cardinality constraints determined by these DTDs are given below:

$$\begin{array}{ll}
 \Psi_{D_1^N}: & \\
 \psi_{teachers}: & |ext(teachers)| = x_{(teacher, teachers)}^1 \qquad |ext(teachers)| = x_{(\tau_t^1, teachers)}^2 \\
 \psi_{\tau_t^1}: & |ext(\tau_t^1)| = x_{(\tau_t^1, \tau_t^1)}^1 + x_{(\tau_t^2, \tau_t^1)}^2 \\
 \psi_{\tau_t^2}: & |ext(\tau_t^2)| = x_{(teacher, \tau_t^2)}^1 \qquad |ext(\tau_t^2)| = x_{(\tau_t^1, \tau_t^2)}^2 \\
 \psi_{teacher}: & |ext(teacher)| = x_{(teach, teacher)}^1 \qquad |ext(teacher)| = x_{(research, teacher)}^2
 \end{array}$$

$$\begin{aligned}
\psi_{\text{teach}}: \quad & |ext(\text{teach})| = x_{(\text{subject}, \text{teach})}^1 & |ext(\text{teach})| &= x_{(\text{subject}, \text{teach})}^2 \\
\psi_{\text{subject}}: \quad & |ext(\text{subject})| = x_{(\mathbf{s}, \text{subject})}^1 \\
\psi_{\text{research}}: \quad & |ext(\text{research})| = x_{(\mathbf{s}, \text{research})}^1 \\
\text{moreover,} \quad & & & \\
& |ext(\text{teachers})| = 1 & |ext(\text{teacher})| &= x_{(\text{teacher}, \text{teachers})}^1 + x_{(\text{teacher}, \tau_t^2)}^1 \\
& |ext(\tau_t^1)| = x_{(\tau_t^1, \text{teachers})}^2 + x_{(\tau_t^1, \tau_t^2)}^2 & |ext(\tau_t^2)| &= x_{(\tau_t^2, \tau_t^1)}^2 \\
& |ext(\tau_\epsilon)| = x_{(\tau_\epsilon, \tau_t^1)}^1 & |ext(\text{teach})| &= x_{(\text{teach}, \text{teacher})}^1 \\
& |ext(\text{subject})| = x_{(\text{subject}, \text{teach})}^1 + x_{(\text{subject}, \text{teach})}^2 \\
& |ext(\text{research})| = x_{(\text{research}, \text{teacher})}^2 & |ext(\mathbf{S})| &= x_{(\mathbf{s}, \text{subject})}^1 + x_{(\mathbf{s}, \text{research})}^1 \\
& \text{all variables} \geq 0.
\end{aligned}$$

For example,  $x_{(\text{teacher}, \text{teachers})}^1$  indicates the number of **teacher** children of all **teachers** nodes, and  $x_{(\text{teacher}, \tau_t^2)}^1$  stands for the number of **teacher** children of nodes labeled  $\tau_t^2$ . The cardinality of  $ext(\text{teacher})$  equals the sum of  $x_{(\text{teacher}, \text{teachers})}^1$  and  $x_{(\text{teacher}, \tau_t^2)}^1$ . Obviously, there is a unique node labeled **teachers**, i.e., the root. Hence we have  $x_{(\text{teacher}, \text{teachers})}^1 = 1$  since the root has a unique **teacher** child. Thus  $|ext(\text{teacher})| = 1 + x_{(\text{teacher}, \tau_t^2)}^1$ .

$$\begin{aligned}
\Psi_{D_2}: \quad & \psi_{db}: \quad |ext(db)| = x_{(foo, db)}^1 \\
& \psi_{foo}: \quad |ext(foo)| = x_{(foo, foo)}^1 \\
& \text{moreover, } |ext(db)| = 1 \quad |ext(foo)| = x_{(foo, db)}^1 + x_{(foo, foo)}^1 \quad \text{all variables} \geq 0.
\end{aligned}$$

It is easy to check that  $\Psi_{D_1^N}$  is consistent, whereas  $\Psi_{D_2^N}$  is not.

We next show that  $\Psi_D$  indeed characterizes the DTD  $D$ .

**LEMMA 4.5.** *Let  $D$  be a simple DTD and  $\Psi_D$  be the set of cardinality constraints determined by  $D$ . Then  $\Psi_D$  is consistent if and only if there is an XML tree  $T$  such that  $T \models D$ . In addition, for each  $\tau \in E$ ,  $|ext(\tau)|$  in  $T$  equals the value of the variable  $|ext(\tau)|$  given by the solution of  $\Psi_D$ .*

*Proof:* First, assume that there is an XML tree  $T$  valid w.r.t.  $D$ . We define an integer solution of  $\Psi_D$  as follows. For each  $\tau \in E \cup \{\mathbf{S}\}$ , let the value of the variable  $|ext(\tau)|$  be the number of  $\tau$  nodes in  $T$ . We proceed to assign integer values (number of certain subelements) to other variables by considering the structure of  $P(\tau)$  for each  $\tau \in E$ . (1) If  $P(\tau) = \tau_1$  for some  $\tau_1 \in E \cup \{\mathbf{S}\}$ , then let the value of the variable  $x_{\tau_1, \tau}^1$  be the number of  $\tau_1$  subelements of all  $\tau$  elements in  $T$ . (2) If  $P(\tau') = (\tau_1, \tau_2)$ , then let the value of the variable  $x_{\tau_1, \tau}^1$  (resp.  $x_{\tau_2, \tau}^2$ ) be the number of the  $\tau_1$  (resp.  $\tau_2$ ) subelements of all  $\tau$  elements. In particular, if  $\tau_1 = \tau_2$ , then

$x_{\tau_1, \tau}^1$  (resp.  $x_{\tau_2, \tau}^2$ ) has the number of the first (resp. second) subelements of all  $\tau$  elements. (3) If  $P(\tau') = (\tau_1 | \tau_2)$ , then let the value of the variable  $x_{\tau_1, \tau}^1$  (resp.  $x_{\tau_2, \tau}^2$ ) be the number of  $\tau_1$  (resp.  $\tau_2$ ) subelements. If  $\tau_1 = \tau_2$ , then  $x_{\tau_1, \tau}^1$  and  $x_{\tau_2, \tau}^2$  may have any value as long as  $|ext(\tau)| = x_{\tau_1, \tau}^1 + x_{\tau_2, \tau}^2$ . We next show that this assignment is an integer solution of  $\Psi_D$ . First, the value of any variable is nonnegative, as it is the number of certain elements (subelements) in  $T$ . Second,  $|ext(r)| = 1$  as  $T$  has a unique root. Third, for each  $\tau \in E$ , by induction on the structure of  $P(\tau)$ , it can be verified that the assignment satisfies  $\psi_\tau$  since  $T \models D$  and  $\psi_\tau$  describes  $P(\tau)$  quantitatively. Finally, the value of the variable  $|ext(\tau)|$  is equal to the sum of all variables of the form  $x_{\tau, \tau'}^i$  ( $i \in [1, 2]$ ) since it counts all the  $\tau$  elements in  $T$  no matter where they are. This can be easily verified by contradiction. Thus the assignment is indeed a solution of  $\Psi_D$ . Note that by the definition of the solution, the value of the variable  $|ext(\tau)|$  given by the solution equals  $|ext(\tau)|$  in  $T$ .

Conversely, assume that  $\Psi_D$  admits an integer solution. Observe that all these variables have nonnegative integer values because of the inequalities in  $\Psi_D$ . We show that there is an XML tree  $T = (V, lab, ele, att, val, root)$  valid w.r.t.  $D$ . To do so, for each  $\tau \in E \cup \{\mathbf{S}\}$ , we create  $|ext(\tau)|$  many distinct nodes and label them with  $\tau$ . We refer to this set of nodes as  $ext(\tau)$ . In addition, for each  $v \in ext(\tau)$  and  $l \in R(\tau)$ , we create a distinct node, referred to as  $v_l$ , and label it with  $l$ . Let

$$\begin{aligned} V &= \bigcup_{\tau \in E \cup \{\mathbf{S}\}} ext(\tau) \cup \bigcup_{\tau \in E} \{v_l \mid v \in ext(\tau), l \in R(\tau)\} \\ lab(v) &= \begin{cases} \tau & \text{if } v \in ext(\tau) \text{ and } \tau \in E \cup \{\mathbf{S}\} \\ l & \text{if } v = v_l \text{ for some } v_l \end{cases} \\ att(v, l) &= \begin{cases} v_l & \text{if } v_l \in V \\ \text{undefined} & \text{otherwise} \end{cases} \\ val(v) &= \begin{cases} \text{empty string} & \text{if } lab(v) \text{ is } \mathbf{S} \text{ or } l, \text{ where } l \in A \\ \text{undefined} & \text{otherwise} \end{cases} \end{aligned}$$

It is easy to verify that these functions are well defined. Let  $root$  be the node labeled  $r$ , which is unique by  $|ext(r)| = 1$  in  $\Psi_D$ . Finally, to define the function  $ele$ , we first mark nodes in  $ext(\tau)$  with variables in  $X_\tau$  so that they can be grouped as subelements of certain elements. For each variable  $x_{\tau, \tau'}^i$  in  $X_\tau$ , we choose  $x_{\tau, \tau'}^i$  many distinct nodes labeled  $\tau$  and mark them with  $x_{\tau, \tau'}^i$ . Note that for each  $\tau \in E \cup \{\mathbf{S}\}$ , every  $\tau$  node in  $V \setminus \{root\}$  can be marked once and only once by  $|ext(\tau)| = \sum_{x_{\tau, \tau'}^i \in X_\tau} x_{\tau, \tau'}^i$  in  $\Psi_D$ . Given these marked elements, starting at  $root$ , for each

$\tau \in E$  and each  $\tau$  node  $v$ , we define  $ele(v)$  as follows. If  $P(\tau)$  is  $\tau_1 \in E \cup \{\mathbf{S}\}$ , then we choose a distinct  $\tau_1$  node  $y$  marked with  $x_{\tau_1, \tau}^1$  and let  $ele(v) = [y]$ . If  $P(\tau) = (\tau_1, \tau_2)$ , then we choose a  $\tau_1$  node  $y_1$  marked with  $x_{\tau_1, \tau}^1$  and a  $\tau_2$  node  $y_2$  marked with  $x_{\tau_2, \tau}^2$ , and let  $ele(v) = [y_1, y_2]$ . If

$P(\tau) = (\tau_1 | \tau_2)$ , then we choose a node  $y$  marked with either  $x_{\tau_1, \tau}^1$  or  $x_{\tau_2, \tau}^2$  and let  $ele(y) = [y]$ . By  $\Psi_D$  constraints, each element or text node in  $V \setminus \{root\}$  can be chosen once and only once as a subelement of some other element. By induction on the structure of  $P(\tau)$ , one can verify that  $T$  defined in this way is indeed an XML tree and  $T \models D$ . Finally, by the definition of  $T$ ,  $|ext(\tau)|$  in  $T$  equals the value of the variable  $|ext(\tau)|$  given by the solution of  $\Psi_D$ .  $\square$

It is straightforward to show that given any simple DTD  $D$ , the set  $\Psi_D$  of cardinality constraints determined by  $D$  can be computed in linear time. As a result, the size of  $\Psi_D$  is linear in  $|D|$ .

**Characterizing DTDs and unary constraints.** To complete our characterization, given a DTD  $D = (E, A, P, R, r)$  and a finite set  $\Sigma$  of  $\mathcal{C}_{K, IC}^{Unary}$  constraints over  $D$ , we define a system  $\Psi(D, \Sigma)$  of integer constraints. The system  $\Psi(D, \Sigma)$ , referred to as *the set of cardinality constraints determined by  $D$  and  $\Sigma$* , is defined to be:

$$\Psi_{D_N} \cup C_\Sigma \cup \{(|ext(\tau)| > 0) \rightarrow (|ext(\tau.l)| > 0) \mid \tau \in E, l \in R(\tau)\},$$

where  $D_N$  is the simplified DTD of  $D$ ,  $\Psi_{D_N}$  and  $C_\Sigma$  are the sets of cardinality constraints determined by  $D_N$  and  $\Sigma$ , respectively. In  $\Psi(D, \Sigma)$  we treat  $|ext(\tau.l)|$  as a variable.

We say that  $\Psi(D, \Sigma)$  is *consistent* if and only if  $\Psi(D, \Sigma)$  admits an integer solution.

For example, recall the DTDs  $D_1$  and  $D_2$ , and the constraint sets  $\Sigma_1$  and  $\Sigma_2$  (the empty set) given in Section 1. It is easy to verify that neither  $\Psi(D_1, \Sigma_1)$  nor  $\Psi(D_2, \Sigma_2)$  is consistent. This is consistent with the observations made in Section 1.

Observe that  $\Psi(D, \Sigma)$  can be partitioned into two sets:  $\Psi(D, \Sigma) = \Psi^l(D, \Sigma) \cup \Psi^e(D, \Sigma)$ , where  $\Psi^l(D, \Sigma)$  consists of linear integer constraints, and  $\Psi^e(D, \Sigma)$  consists of constraints of the form  $(|ext(\tau)| > 0 \rightarrow |ext(\tau.l)| > 0)$ , which are to ensure that every  $\tau$  element has an  $l$  attribute. Note that  $|ext(\tau.l)| \leq |ext(\tau)|$  is already in  $C_\Sigma$ .

It is easy to verify that  $\Psi(D, \Sigma)$  can be computed in linear time in  $|D|$  and  $|\Sigma|$ , and thus its size is also linear in  $|D|$  and  $|\Sigma|$ .

We next show that  $\Psi(D, \Sigma)$  indeed characterizes  $D$  and  $\Sigma$ .

**LEMMA 4.6.** *Let  $D$  be a DTD,  $\Sigma$  be a finite set of  $\mathcal{C}_{K, IC}^{Unary}$  constraints over  $D$ , and  $\Psi(D, \Sigma)$  be the set of cardinality constraints determined by  $D$  and  $\Sigma$ . Then  $\Psi(D, \Sigma)$  is consistent if and only if there exists an XML tree  $T$  such that  $T \models D$  and  $T \models \Sigma$ .*

*Proof:* Let  $D_N$  be the simplified DTD of  $D$ . By Lemma 4.3, it suffices to show that  $\Psi(D, \Sigma)$  is consistent if and only if there is an XML tree  $T$  such that  $T \models D_N$  and  $T \models \Sigma$ .

Suppose that there exists an XML tree  $T$  such that  $T \models D_N$  and  $T \models \Sigma$ . We show that



$\Psi(D, \Sigma)$  admits an integer solution. By Lemma 4.4, we have  $T \models C_\Sigma$ , where  $C_\Sigma$  is the set of cardinality constraints determined by  $\Sigma$ . By Lemma 4.5, one can define an integer solution of  $\Psi_{D_N}$ . The assignment assures that for each  $\tau \in E$ , the value of the variable  $|ext(\tau)|$  equals the number of all the  $\tau$  nodes in  $T$ . We extend the assignment as follows: for each  $\tau \in E$  and  $l \in R(\tau)$ , let the value of the variable  $|ext(\tau.l)|$  be the number of distinct  $l$  attribute values of all the  $\tau$  nodes in  $T$ . Thus by  $T \models C_\Sigma$ , this extended assignment satisfies  $C_\Sigma$ . In addition, if  $|ext(\tau)| > 0$  then  $|ext(\tau.l)| > 0$  as every  $\tau$  element in  $T$  has an  $l$  attribute. Hence the assignment is indeed a solution to  $\Psi(D, \Sigma)$ . Thus  $\Psi(D, \Sigma)$  is consistent.

Conversely, suppose that  $\Psi(D, \Sigma)$  admits an integer solution. We show that there is an XML tree  $T$  such that  $T \models D_N$  and  $T \models \Sigma$ . Observe that an integer solution to  $\Psi(D, \Sigma)$  is also a solution to  $\Psi_{D_N}$ . Thus by Lemma 4.5, there is  $T' = (V, lab, ele, att, val, root)$  such that  $T' \models D_N$ . Moreover, for each  $\tau \in E$ ,  $|ext(\tau)|$  in  $T'$  is equal to the value of the variable  $|ext(\tau)|$  given by the assignment. We construct another XML tree  $T''$  by modifying the definition of the function  $val$  of  $T'$  such that for each  $\tau \in E$  and  $l \in R(\tau)$ ,  $|ext(\tau.l)|$  in  $T''$  equals the value assigned to the variable  $|ext(\tau.l)|$  by the assignment. This is possible since  $|ext(\tau.l)| \leq |ext(\tau)|$  is in  $C_\Sigma$ , and the assignment is also a solution to  $C_\Sigma$ . Moreover, by  $(|ext(\tau)| > 0 \rightarrow |ext(\tau.l)| > 0)$  in  $\Psi(D, \Sigma)$ , every  $\tau$  element in  $T''$  can have an  $l$  attribute. It is straightforward to verify that  $T'' \models C_\Sigma$  and  $T'' \models D_N$ . Hence by Lemma 4.4, there exists an XML tree  $T$  such that  $T \models D_N$  and  $T \models \Sigma$ .  $\square$

Given these lemmas, we proceed to prove Theorem 4.1.

*Proof of Theorem 4.1 (continued):* We encode an instance  $(D, \Sigma)$  of the consistency problem for  $\mathcal{C}_{K,FK}^{Unary}$  as an instance of LIP. By Lemma 4.6, it suffices to encode  $\Psi(D, \Sigma)$  as an instance of LIP. Recall that  $\Psi(D, \Sigma)$  can be partitioned into two sets:  $\Psi^l(D, \Sigma)$  of linear integer constraints, and  $\Psi^c(D, \Sigma)$  of constraints of the form  $(x > 0 \rightarrow y > 0)$ . We first encode  $\Psi(D, \Sigma)$  with a set of linear integer constraints. Let  $S$  be the set of all the pairs  $(x, y)$  for each constraint  $(x > 0 \rightarrow y > 0)$  in  $\Psi^c(D, \Sigma)$ . For each subset  $X$  of  $S$ , we define  $\Psi_X$  to be

$$\Psi^l(D, \Sigma) \cup \{x = 0, y = 0 \mid (x, y) \in X\} \cup \{x \geq 1, y \geq 1 \mid (x, y) \in S \setminus X\}.$$

It is easy to see that  $\Psi(D, \Sigma)$  admits an integer solution if and only if there is some  $\Psi_X$  that has an integer solution. Observe that  $\Psi_X$  can be represented as an instance of LIP since an equality  $F_1 = F_2$  is equivalent to inequalities  $F_1 \geq F_2$  and  $F_2 \geq F_1$ . In addition, for all variables  $x$  in  $\Psi(D, \Sigma)$ , we have  $x \geq 0$  in  $\Psi(D, \Sigma)$ . Thus any solution of  $\Psi_X$  is nonnegative. Hence we can apply the result of [32] here, which says that if  $\Psi_X$  has an integer solution, then it has one in which the values of all variables are no larger than  $n(ma)^{2m+1}$ , where  $a$  is the largest absolute value of the constants in  $\Psi_X$ . In other words,  $\Psi_X$  has an integer solution in which the value of each variable has a length in binary of at most  $1 + \lceil \log n + (2m + 1) \cdot \log(ma) \rceil$  many bits, and the bounds on solutions for all  $\Psi_X$ 's are the same. Let  $c$  be a number that in binary

notation has  $1 + \lceil \log n + (2m + 1) \cdot \log(ma) \rceil$  many 1's. Observe that  $c$  can be computed in  $O(s \log s)$  time. Thus we define a new system  $\Phi$  of linear integer constraints that is the same as  $\Psi^l(D, \Sigma)$  except it also includes  $cy \geq x$  for all  $(x > 0) \rightarrow (y > 0)$  in  $\Psi^c(D, \Sigma)$ . It is easy to verify that  $\Psi(D, \Sigma)$  has an integer solution iff  $\Phi$  has an integer solution. Indeed, if  $\Psi(D, \Sigma)$  has an integer solution then it has one bounded by  $c$ . Thus the solution satisfies  $cy \geq x$ , i.e., it is an integer solution to  $\Phi$ . Conversely, if  $\Phi$  has an integer solution, then it is also an integer solution of  $\Psi^l(D, \Sigma)$  and moreover, if  $x > 0$  then  $y > 0$  by  $cy \geq x$  in  $\Phi$ ; that is, it is an integer solution to  $\Psi(D, \Sigma)$ . As  $\Phi$  can be represented as an instance of LIP, we can define an matrix  $A_\Psi$  and a vector  $\vec{b}_\Psi$  of integers such that  $\Psi(D, \Sigma)$  has an integer solution if and only if  $A_\Psi \vec{x} \geq \vec{b}_\Psi$  has an integer solution. Recall that  $\Psi(D, \Sigma)$  can be computed in linear time and its size, denoted by  $s$ , is linear in  $|D|$  and  $|\Sigma|$ . Thus the instance of LIP can be computed in  $O(s^2 \cdot \log s)$  time in  $|D|$  and  $|\Sigma|$ .

This completes the proof of Theorem 4.1.  $\square$

The encoding is not only interesting in its own right, but also useful in the consistency analyses of  $\mathcal{C}_{K,FK}^{Unary}$  and  $\mathcal{C}_{K^\neg,IC}^{Unary}$  constraints, as well as in resolving a special case of  $\mathcal{C}_{K,FK}^{Unary}$  constraint implication.

#### 4.2 $\mathcal{C}_{K,FK}^{Unary}$ and $\mathcal{C}_{K^\neg,IC}^{Unary}$ constraints

We next establish the precise complexity bound on the consistency problem for unary keys and foreign keys:

**THEOREM 4.7.** *The consistency problem for  $\mathcal{C}_{K,FK}^{Unary}$  constraints is NP-complete.*

*Proof:* Corollary 4.2 has shown that the problem is in NP. We show that it is NP-hard by reduction from a variant of LIP, namely,

$$A \vec{x} = \vec{b},$$

where for all  $i \in [1, m]$ ,  $j \in [1, n]$ ,  $a_{ij}$  coefficients are in  $\{0, 1\}$ , all  $b_i$  elements are 1, and all  $x_j$  components are binary, i.e., in  $\{0, 1\}$ . It is known that the variant is also NP-complete [22].

Given such an instance  $A \vec{x} = \vec{b}$ , we define a DTD  $D$  and a set  $\Sigma$  of  $\mathcal{C}_{K,FK}^{Unary}$  constraints over  $D$  such that there is an XML tree valid w.r.t.  $D$  and satisfying  $\Sigma$  if and only if  $A \vec{x} = \vec{b}$  admits a binary solution. For  $i \in [1, m]$ , we use  $F_i$  to denote  $\sum_{j \in [1, n]} a_{ij} x_j$ . We define  $D$  to be

$(E, A, P, R, r)$ , where

$$E = \{r\} \cup \{F_i \mid i \in [1, m]\} \cup \{b_i \mid i \in [1, m]\} \cup \{VF_i \mid i \in [1, m]\}$$

$$\begin{aligned}
& \cup \{X_{ij} \mid i \in [1, m], j \in [1, n]\} \cup \{Z_{ij} \mid i \in [1, m], j \in [1, n]\} \\
A = & \{v\} \cup \{A_{ij} \mid i \in [1, m], j \in [1, n]\} \\
P(r) = & F_1, \dots, F_m, b_1, \dots, b_m \\
P(F_i) = & X_{ij_1}, \dots, X_{ij_l} \quad \text{for } i \in [1, m], \text{ where } X_{ij_1}, \dots, X_{ij_l} \text{ is a sub-list of } X_{i1}, \dots, X_{im} \\
& \text{such that } X_{ij} \text{ is in } P(F_i) \text{ iff } a_{ij} \text{ in } A \text{ is } 1 \\
P(X_{ij}) = & Z_{ij} \mid \epsilon \quad \text{for } i \in [1, m] \text{ and } j \in [1, n] \\
P(Z_{ij}) = & VF_i \quad \text{for } i \in [1, m] \text{ and } j \in [1, n] \\
P(VF_i) = & P(b_i) = \epsilon \quad \text{for } i \in [1, m] \\
R(Z_{ij}) = & \{A_{ij}\} \quad \text{for } i \in [1, m] \text{ and } j \in [1, n] \\
R(VF_i) = & R(b_i) = \{v\} \quad \text{for } i \in [1, m] \\
R(r) = & R(F_i) = R(X_{ij}) = \emptyset
\end{aligned}$$

An XML tree valid w.r.t.  $D$  has the form shown in Figure 4. Intuitively,  $X_{ij}$  encodes  $x_j$  in  $F_i$ , and  $Z_{ij}$  encodes the value of  $X_{ij}$ :  $X_{ij}$  has value 1 if and only if  $X_{ij}$  has a  $Z_{ij}$  child. The element type  $VF_i$  is to code the value of  $F_i$ . Observe that  $A\vec{x} = \vec{b}$  has a solution if and only if for each row  $i \in [1, m]$  there is exactly one column  $j \in [1, n]$  such that  $a_{ij} = 1$  and  $x_j = 1$ . In the XML tree  $T$  representing the instance, this means that for every  $i$  there is exactly one  $X_{ij}$  element with a  $Z_{ij}$  child. This is achieved by restricting  $F_i$  to have a unique  $VF_i$  descendant, and thus to have value 1, by means of the attribute  $v$  of  $VF_i$  and constraints. More specifically, we include the following in the set  $\Sigma$ :

$$VF_i.v \rightarrow VF_i, \quad b_i.v \rightarrow b_i, \quad VF_i.v \subseteq b_i.v, \quad b_i.v \subseteq VF_i.v.$$

These ensure that  $F_i = b_i = 1$  as  $T$  has a unique  $b_i$  node. In addition, to ensure that all occurrences of  $x_j$  have the same value, the following are in  $\Sigma$ : for  $j \in [1, n]$  and  $i, l \in [1, m]$ ,

$$Z_{ij}.A_{ij} \rightarrow Z_{ij}, \quad Z_{ij}.A_{ij} \subseteq Z_{lj}.A_{lj}.$$

These assert that  $X_{ij}$  has value 1 if and only if  $X_{lj}$  equals 1. It is easy to see that the encoding can be done in PTIME in  $m$  and  $n$ . Moreover,  $A\vec{x} = \vec{b}$  admits a binary solution if and only if  $D$  has a valid XML tree satisfying  $\Sigma$ . Thus this is indeed a PTIME reduction from the variant of LIP.  $\square$

Recall that in relational databases, it is common to consider primary keys. That is, for each relation one can specify at most one key, namely, the primary key of the relation. In the XML setting, the *primary key restriction* requires that for each element type one can specify at most one key. This is the case for “keys” specified with ID attributes, since in a DTD, at most one ID attribute can be specified for each element type. Under the primary key restriction, the consistency problem for a class  $\mathcal{C}$  of XML constraints is to determine, given any DTD  $D$  and

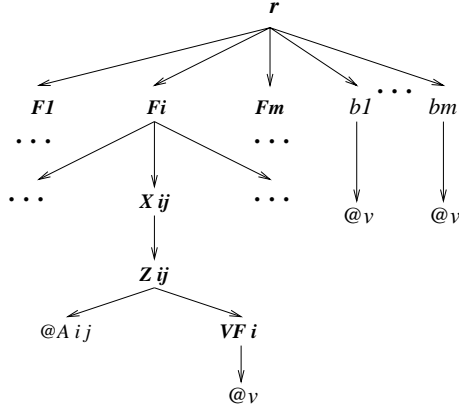


Fig. 4. A tree used in the proof of Theorem 4.7

finite set  $\Sigma$  of  $\mathcal{C}$  constraints in which there is at most one key for each element type (given either as keys or as part of foreign keys), whether there is an XML tree valid w.r.t.  $D$  and satisfying  $\Sigma$ ; similarly for implication.

One might think that the primary key restriction would simplify the consistency analysis of  $\mathcal{C}_{K,FK}^{Unary}$  constraints. However, it is not the case.

**COROLLARY 4.8.** *Under the primary key restriction, the consistency problem for  $\mathcal{C}_{K,FK}^{Unary}$  remains NP-complete.*

*Proof:* The reduction from LIP given in the proof of Theorem 4.7 defines at most one key for each element type.  $\square$

A mild generalization of the encoding above can establish the complexity of the consistency problem for  $\mathcal{C}_{K^-,IC}^{Unary}$ , the class of unary keys, inclusion constraints and negations of keys. As we shall see shortly, the result for  $\mathcal{C}_{K^-,IC}^{Unary}$  helps us study implication of  $\mathcal{C}_{K,FK}^{Unary}$  constraints.

**COROLLARY 4.9.** *The consistency problem for  $\mathcal{C}_{K^-,IC}^{Unary}$  constraints is NP-complete.*

*Proof:* Since  $\mathcal{C}_{K,FK}^{Unary}$  is a sub-language of  $\mathcal{C}_{K^-,IC}^{Unary}$ , from Theorem 4.7 follows immediately that the consistency problem for  $\mathcal{C}_{K^-,IC}^{Unary}$  is NP-hard. We next show that the problem remains in NP. Let  $D$  be a DTD and  $\Sigma$  be a set of  $\mathcal{C}_{K^-,IC}^{Unary}$  constraints over  $D$ . We write  $\Sigma$  as  $\Sigma_1 \cup \Sigma_2$ , where  $\Sigma_1$  is a set of unary keys and unary inclusion constraints over  $D$ , and  $\Sigma_2$  is a set of negations of unary keys over  $D$ . Let  $\Psi(D, \Sigma_1)$  be the system of linear inequalities determined by  $D$  and  $\Sigma_1$ , as defined in the proof of Theorem 4.1. It admits an integer solution iff there exists an XML

tree  $T$  such that  $T \models \Sigma_1$  and  $T \models D$ . We define another system of linear inequalities, denoted by  $\Psi(D, \Sigma)$  and referred to as *the system determined by  $D$  and  $\Sigma$* , to be

$$\Psi(D, \Sigma) = \Psi(D, \Sigma_1) \cup \{ |ext(\tau.l)| < |ext(\tau)| \mid \neg(\tau.l \rightarrow \tau) \in \Sigma_2 \}.$$

As  $\Psi(D, \Sigma)$  can be computed in PTIME, it suffices to show the following claim.

*Claim:* There is an XML tree  $T$  such that  $T \models \Sigma$  and  $T \models D$  iff  $\Psi(D, \Sigma)$  has an integer solution.

For if it holds, then the problem is in NP by reduction to LIP as in the proof of Theorem 4.1.

We show the claim as follows. Assume that there exists a tree  $T$  such that  $T \models \Sigma$  and  $T \models D$ . Since  $T \models \Sigma_1$ , by Lemmas 4.5 and 4.6 and Theorem 4.1, it can be verified that there is an integer solution to  $\Psi(D, \Sigma_1)$ , the system of linear inequalities determined by  $D$  and  $\Sigma_1$ , such that the values of the variables  $|ext(\tau)|$  and  $|ext(\tau.l)|$  in  $\Psi(D, \Sigma_1)$  given by the solution are the cardinalities  $|ext(\tau)|$  and  $|ext(\tau.l)|$  in  $T$ . Note that for all element type  $\tau$  and attribute  $l$  of  $\tau$  in  $D$ ,  $|ext(\tau)|$  and  $|ext(\tau.l)|$  are variables in  $\Psi(D, \Sigma_1)$ . Thus for each  $\tau.l \not\rightarrow \tau$ , the solution also assigns values to  $|ext(\tau)|$  and  $|ext(\tau.l)|$ . We claim that it is also a solution to  $\Psi(D, \Sigma)$ . To see this, observe that it is always true that  $|ext(\tau)| \geq |ext(\tau.l)|$  in  $T$  since every  $\tau$  element in  $T$  contributes at most one distinct  $\tau.l$  value. Thus by  $T \models \Sigma_2$ , there must be two distinct  $\tau$  elements  $d_1$  and  $d_2$  in  $T$  such that  $d_1.l = d_2.l$ . Thus  $|ext(\tau)| > |ext(\tau.l)|$ . Therefore, all inequalities in  $\Psi(D, \Sigma)$  are satisfied by the solution.

Conversely, assume that  $\Psi(D, \Sigma)$  has an integer solution. Since it is also a solution to  $\Psi(D, \Sigma_1)$ , again by Lemma 4.5 and 4.6 and Theorem 4.1, it can be verified that there is a tree  $T$  such that  $T \models D$ ,  $T \models \Sigma_1$  and moreover, the cardinalities  $|ext(\tau)|$  and  $|ext(\tau.l)|$  in  $T$  are the values of the variables  $|ext(\tau)|$  and  $|ext(\tau.l)|$  in  $\Psi(D, \Sigma_1)$  given by the solution. We claim that  $T \models \Sigma$ . Indeed, for any  $\tau.l \not\rightarrow \tau$  in  $\Sigma_2$ , we have  $|ext(\tau)| > |ext(\tau.l)|$  in  $T$ . Thus there must be two distinct  $\tau$  elements  $d_1$  and  $d_2$  in  $T$  such that  $d_1.l = d_2.l$ . That is,  $T \models \tau.l \not\rightarrow \tau$ . Hence  $T \models D$  and  $T \models \Sigma$ .  $\square$

It should be mentioned that the problem remains NP-hard under the primary key restriction. This can be verified along the same lines as the proof of Corollary 4.8.

Corollary 4.9 also tells us the complexity of a special case of the implication problem for  $\mathcal{C}_{K,FK}^{Unary}$ , referred to as *implication problem for unary keys by  $\mathcal{C}_{K,FK}^{Unary}$  constraints*:

**THEOREM 4.10.** *The following is coNP-complete, even under the primary key restriction: given any DTD  $D$ , any set  $\Sigma$  of  $\mathcal{C}_{K,FK}^{Unary}$  constraints and any unary key  $\varphi$  over  $D$ , whether  $(D, \Sigma) \vdash \varphi$ .*

*Proof:* Observe that  $(D, \Sigma) \vdash \varphi$  iff  $\Sigma \cup \{\neg\varphi\}$  and  $D$  are not consistent, i.e., there exists no XML tree  $T$  such that  $T \models D$ ,  $T \models \Sigma$  and  $T \models \neg\varphi$ . Since  $\Sigma \cup \{\neg\varphi\}$  is a set of  $\mathcal{C}_{K^-,IC}^{Unary}$  constraints, the implication problem for unary keys by  $\mathcal{C}_{K,FK}^{Unary}$  constraints is in coNP by Corollary 4.9. To see that the problem is coNP-hard, recall the encoding given in the proof of Lemma 3.3. If the set  $\Sigma$  of constraints given is a set of  $\mathcal{C}_{K,FK}^{Unary}$  constraints, then that encoding also serves as a reduction from the consistency problem for  $\mathcal{C}_{K,FK}^{Unary}$  to the complement of  $(D, \Sigma) \vdash \varphi$ . Thus from Theorem 4.1 follows that the implication problem for unary keys by  $\mathcal{C}_{K,FK}^{Unary}$  constraints is coNP-hard. Observe that the reduction in the proof of Lemma 3.3 defines at most one key for each element type. Thus given a set  $\Sigma$  of constraints, if  $\Sigma$  satisfies the primary key restriction, then so does the set of all constraints used in the reduction. Hence it remains coNP-hard even under the primary key restriction.  $\square$

Finally, we identify some PTIME decidable cases of the consistency and implication problems. First, these problems for unary keys only are decidable in linear time, by Theorem 3.5. We next show that given a fixed DTD  $D$ , the consistency and implication analyses become simpler. The motivation for considering a fixed DTD is because in practice, one often defines the DTD of a specification at one time, but writes constraints in stages: constraints are added incrementally when new requirements are discovered.

**COROLLARY 4.11.** *For a fixed DTD, the following problems are decidable in PTIME:*

- The consistency problems for  $\mathcal{C}_{K,FK}^{Unary}$  and  $\mathcal{C}_{K^-,IC}^{Unary}$ .
- Implication of unary keys by  $\mathcal{C}_{K,FK}^{Unary}$  constraints.

*Proof:* By Theorems 4.1, 4.10 and Corollary 4.9, an instance  $(D, \Sigma)$  of these problems can be encoded as a system  $\Phi$  of linear integer constraints. That is, these problems can be reduced to checking whether  $\Phi$  admits an integer solution. The system  $\Phi$  consists of constraints of  $C_\Sigma$  (derived from  $\Sigma$ ) and  $\Psi_{D_N}$  (derived from the simplified DTD  $D_N$  of  $D$ ), and can be computed in PTIME in  $|D|$ . Given a fixed DTD  $D$ , the number of variables in  $C_\Sigma$  is bounded by the size of  $D$  ( $O(|D|^2)$ ), and the number of variables in  $\Psi_{D_N}$  is also fixed. Thus the number of variables in  $\Phi$  is bounded. It is known that when the number of variables in a system of linear integer constraints is bounded, checking whether the system admits an integer solution can be done in PTIME [29]. Putting these together, we have Corollary 4.11.  $\square$

## 5. UNARY KEYS, INCLUSION CONSTRAINTS AND NEGATIONS

In Section 4, we have shown that the consistency problem for unary keys and foreign keys is NP-complete. In this section, we extend the result by showing that the problem remains in NP when negations of these unary constraints are allowed. That is, the problem is NP-complete

for  $\mathcal{C}_{K^-,IC^-}^{Unary}$ , the class of unary keys, inclusion constraints and their negations. This helps us settle the implication problems for  $\mathcal{C}_{K,FK}^{Unary}$  and the more general  $\mathcal{C}_{K,IC}^{Unary}$ , the class of unary keys and foreign keys, and the class of unary keys and inclusion constraints, respectively. This is one of the reasons that we are interested in the consistency problem for  $\mathcal{C}_{K^-,IC^-}^{Unary}$ .

**THEOREM 5.1.** *The consistency problem for  $\mathcal{C}_{K^-,IC^-}^{Unary}$  is NP-complete.*

While this theorem subsumes Theorem 4.7, the reduction is quite different from the nice encoding with instances of LIP that we used for  $\mathcal{C}_{K,FK}^{Unary}$ . In fact, while typically NP-complete problems are easily shown to be in NP, and only the reduction from a known NP-complete problem is difficult, for the consistency problem for  $\mathcal{C}_{K^-,IC^-}^{Unary}$ , the opposite is the case, and the proof of membership in NP is a little involved (even assuming the encoding of keys and inclusion constraints by instances of LIP given in the previous section). We cannot reduce the problem directly to LIP as before, because there is no direct connection between  $\tau_i.l_i \not\subseteq \tau_j.l_j$  and the cardinalities  $|ext(\tau_i)|$ ,  $|ext(\tau_j)|$ ,  $|ext(\tau_i.l_i)|$  and  $|ext(\tau_j.l_j)|$  in an XML tree.

*Proof:* We develop an NP algorithm for determining the consistency of  $\mathcal{C}_{K^-,IC^-}^{Unary}$  constraints. The algorithm takes advantage of another encoding of  $\mathcal{C}_{K^-,IC^-}^{Unary}$  constraints with linear integer constraints, which characterizes a set interpretation of unary inclusion constraints and their negations. Let  $D$  be a DTD and  $\Sigma$  be a set of  $\mathcal{C}_{K^-,IC^-}^{Unary}$  constraints over  $D$ . We partition  $\Sigma$  into  $\Sigma_1$  and  $\Sigma_2$ , where  $\Sigma_1$  is a set of  $\mathcal{C}_{K^-,IC}^{Unary}$  constraints, and  $\Sigma_2$  consists of negations of unary inclusion constraints over  $D$ . Let  $\Psi(D, \Sigma_1)$  be the system of linear inequalities determined by  $D$  and  $\Sigma_1$ , as described in the proof of Corollary 4.9. Let  $l_1, \dots, l_n$  be an enumeration of all attributes in  $D$ . Without loss of generality, assume that  $l_i$  is an attribute of element type  $\tau_i$  (note that  $\tau_i$ 's need not be distinct). Let  $\mathbf{U} = (u_{ij})_{i,j=1}^n$  and  $\mathbf{V} = (v_{ij})_{i,j=1}^n$  be two matrices whose elements are nonnegative integers. We say that they admit a *set representation* if there is a family of finite sets  $A_1, \dots, A_n$  such that

$$u_{ij} = |A_i \cap A_j|, \quad v_{ij} = |A_i \setminus A_j|.$$

We extend  $\Psi(D, \Sigma_1)$  with new variables  $u_{ij}, v_{ij}$ , and equalities:

- $|ext(\tau_i.l_i)| = u_{ii} = u_{ij} + v_{ij}$  for all  $i, j \in [1, n]$ ;
- $v_{ij} = 0$  for all  $\tau_i.l_i \subseteq \tau_j.l_j$  in  $\Sigma_1$ , and moreover,  $v_{ii} = 0$ ;
- $v_{ij} > 0$  for all  $\tau_i.l_i \not\subseteq \tau_j.l_j$  in  $\Sigma_2$ .

Let us denote the new system by  $\Psi(D, \Sigma)$  and refer to it as *the system determined by  $D$  and  $\Sigma$* . Observe that  $\Psi(D, \Sigma)$  can be simply converted to a system of linear inequalities (by treating an equality as two inequalities).

The intended interpretation for the variable  $u_{ij}$  is  $|ext(\tau_i.l_i) \cap ext(\tau_j.l_j)|$ , and  $|ext(\tau_i.l_i) \setminus ext(\tau_j.l_j)|$  for  $v_{ij}$ . Thus  $v_{ij} > 0$  in  $\Psi(D, \Sigma)$  says that  $ext(\tau_i.l_i) \not\subseteq ext(\tau_j.l_j)$  for all  $\tau_i.l_i \not\subseteq \tau_j.l_j$  in  $\Sigma_2$ .

The lemma below reveals the connection between the encoding and the consistency problem we are investigating.

**LEMMA 5.2.** *The linear system  $\Psi(D, \Sigma)$  determined by DTD  $D$  and constraints  $\Sigma$  has an integer solution with  $\mathbf{U}, \mathbf{V}$  having a set representation if and only if there is an XML tree  $T$  such that  $T \models D$  and  $T \models \Sigma$ .*

*Proof:* Let  $D$  be a DTD,  $\Sigma_1$  be a set of  $\mathcal{C}_{K^-, IC}^{Unary}$  constraints over  $D$ ,  $\Sigma_2$  be a set of negations of unary inclusion constraints over  $D$ ,  $\Sigma = \Sigma_1 \cup \Sigma_2$ , and  $\Psi(D, \Sigma)$  be the system of linear inequalities determined by  $D$  and  $\Sigma$  as described above. We show that  $\Psi(D, \Sigma)$  has an integer solution with  $\mathbf{U}, \mathbf{V}$  having a set representation iff there is an XML tree  $T$  such that  $T \models \Sigma$  and  $T \models D$ .

Assume that there exists an XML tree  $T$  such that  $T \models \Sigma$  and  $T \models D$ . Since  $T \models \Sigma_1$ , as in the proof of Corollary 4.9 we can define an integer solution to  $\Psi(D, \Sigma_1)$ , the system of linear inequalities determined by  $D$  and  $\Sigma_1$ . We extend the solution as follows: let  $u_{ij}$  be  $|ext(\tau_i.l_i) \cap ext(\tau_j.l_j)|$ , and  $v_{ij}$  be  $|ext(\tau_i.l_i) \setminus ext(\tau_j.l_j)|$ . It is easy to verify that this is indeed a solution to  $\Psi(D, \Sigma)$  with  $\mathbf{U}, \mathbf{V}$  having a set representation.

Conversely, assume that  $\Psi(D, \Sigma)$  has an integer solution with  $\mathbf{U}, \mathbf{V}$  having a set representation. Then there are finite sets  $A_1, \dots, A_n$  such that

$$u_{ij} = |A_i \cap A_j|, \quad v_{ij} = |A_i \setminus A_j|.$$

Again as in the proof of Corollary 4.9, we create a tree  $T$  such that  $T \models \Sigma_1$  and  $T \models D$ . In addition, we define the *val* function in  $T$  such that  $ext(\tau_i.l_i) = A_i$  for  $i \in [1, n]$ . This is possible since  $|ext(\tau_i.l_i)| = u_{ii} = u_{ij} + v_{ij}$  is in  $\Psi(D, \Sigma)$  for all  $i, j \in [1, n]$ . Because  $v_{ij} > 0$  is in  $\Psi(D, \Sigma)$  for all  $\tau_i.l_i \not\subseteq \tau_j.l_j$  in  $\Sigma_2$ , we have  $|ext(\tau_i.l_i) \setminus ext(\tau_j.l_j)| > 0$ . That is,  $T \models \tau_i.l_i \not\subseteq \tau_j.l_j$ . Thus  $T \models \Sigma_2$ . This completes the proof of the lemma.  $\square$

It remains to show that one can check in NP whether the system  $\Psi(D, \Sigma)$  has an integer solution with  $\mathbf{U}, \mathbf{V}$  having a set representation. We start with a lemma.

**LEMMA 5.3.** *Given  $\Psi(D, \Sigma)$ , one can compute, in polynomial time, a number  $M$  such that  $\Psi(D, \Sigma)$  has an integer solution with  $\mathbf{U}, \mathbf{V}$  having a set representation if and only if it admits such a solution with all variables being bounded by  $M$ .*

*Proof:* To prove the lemma, we need to extend  $\Psi(D, \Sigma)$ . Let  $\Theta$  be the set of functions  $\theta : \{1, \dots, n\} \rightarrow \{0, 1\}$  which are not identically 0, where  $n$  is the number of attributes in  $D$ . For



every  $\theta$ , we introduce a new variable  $z_\theta$  (note that the number of variables is now exponential in the size of the problem). The intended interpretation of  $z_\theta$  is the cardinality of

$$\bigcap_{i:\theta(i)=1} ext(\tau_i.l_i) \setminus \bigcup_{j:\theta(j)=0} ext(\tau_j.l_j).$$

We now extend  $\Psi(D, \Sigma)$  to  $\Psi'(D, \Sigma)$  by adding the following equalities:

$$u_{ij} = \sum_{\theta:\theta(i)=\theta(j)=1} z_\theta, \quad v_{ij} = \sum_{\theta:\theta(i)=1, \theta(j)=0} z_\theta.$$

Clearly,  $\Psi(D, \Sigma)$  has an integer solution with  $\mathbf{U}, \mathbf{V}$  having a set representation iff  $\Psi'(D, \Sigma)$  has an integer solution, as the variables  $z_\theta$  describe all possible intersections of  $ext(\tau_i.l_i)$  and their complements, and the equalities above show how to reconstruct  $u_{ij}$  and  $v_{ij}$  from them. We thus must show that if  $\Psi'(D, \Sigma)$  has an integer solution then it must have one with a bound on  $u_{ij}, v_{ij}$ , which is polynomial (in terms of the size of  $\Psi(D, \Sigma)$ ). For that, recall [32] that if a system of  $k$  linear inequalities with  $l$  variables and all coefficients at most  $c$  has an integer solution, then it has an integer solution in which none of the variables exceeds  $l(ck)^{2k+1}$ . Thus,  $M$  can be taken to be a number that in binary notation has  $1 + \lceil \log l + (2k + 1) \cdot \log(ck) \rceil$  many 1's. Note that the number of variables,  $l$ , of  $\Psi'(D, \Sigma)$  is at most exponential in the size of  $\Psi(D, \Sigma)$ , and the number of equalities,  $k$ , is at most polynomial. This shows that  $M$  can be found in polynomial time, and thus proves the lemma.  $\square$

Given Lemmas 5.2 and 5.3, let us go back to the proof of that consistency analysis of  $\Sigma$  over  $D$  is in NP. We present an NP algorithm for determining the consistency of  $\Sigma$  over  $D$ . Our nondeterministic machine computes  $M$  given by Lemma 5.3, and then guesses a solution with all the components bounded by  $M$ . It then tests if the  $\mathbf{U}, \mathbf{V}$  part has a set representation. To do so, we transform  $\mathbf{U}, \mathbf{V}$ , in polynomial time, into another matrix  $\mathbf{W}$ , and then run a nondeterministic polynomial time machine on  $\mathbf{W}$ . If it returns ‘yes’, then  $\mathbf{U}, \mathbf{V}$  have a set representation, and thus by Lemma 5.2 the answer to whether  $\Sigma$  is consistent over  $D$  is ‘yes’.

Let  $K = M \cdot n$ , where  $n$  is the number of all attributes in  $D$ . We now define the matrix  $\mathbf{W}$ . It is a  $2n \times 2n$  matrix, with

$$w_{ij} = \begin{cases} u_{ij} & \text{if } i, j \leq n \\ v_{i, j-n} & \text{if } i \leq n, j > n \\ v_{i-n, j} & \text{if } i > n, j \leq n \\ K - u_{i-n, j-n} - v_{i-n, j-n} - v_{j-n, i-n} & \text{if } i, j > n \end{cases}$$

Recall the INTERSECTION PATTERN problem: Given an  $m \times m$  matrix  $\mathbf{A}$ , are there sets  $Y_1, \dots, Y_m$  such that  $a_{ij} = |Y_i \cap Y_j|$ ? This problem is known to be NP-complete (see, e.g., [22]).

We now show the following: The INTERSECTION PATTERN problem returns ‘yes’ on input  $\mathbf{W}$  iff  $\mathbf{U}, \mathbf{V}$  have a set representation.

First, assume  $\mathbf{U}, \mathbf{V}$  have a set representation. That is, there are finite sets  $A_1, \dots, A_n$  such that

$$u_{ij} = |A_i \cap A_j|, \quad v_{ij} = |A_i \setminus A_j|.$$

By the assumption, all entries in  $\mathbf{U}, \mathbf{V}$  are bounded by  $M$ , and hence we may assume that all sets in the representation are subsets of a set  $U$  of cardinality  $K$ . Let  $m = 2n$  and define  $Y_i$  to be  $A_i$  for  $i \leq n$ , and  $U \setminus A_{i-n}$  for  $i > n$ . Then  $\mathbf{W}$  is the intersection pattern for this family of sets, and thus the INTERSECTION PATTERN problem returns ‘yes’ on  $\mathbf{W}$ .

Next, assume that the INTERSECTION PATTERN returns ‘yes’ on  $\mathbf{W}$ , so we have a family of sets  $Y_1, \dots, Y_{2n}$  for which  $\mathbf{W}$  is the intersection pattern. Let  $U$  be the union of all  $Y_j$ ’s. We show  $Y_{n+i} = U \setminus Y_i$  for all  $i \leq n$ . We have  $w_{i,n+i} = v_{ii} = 0$ , and thus  $Y_{n+i} \subseteq U \setminus Y_i$ . Moreover, we have  $|Y_i \cup Y_{n+i}| = w_{ii} + w_{n+i,n+i} = K$ . We next show that for every  $i, j \leq n$  it is the case that  $Y_i \cup Y_{n+i} = Y_j \cup Y_{n+j}$  (and thus equals  $U$ ). Note that both  $Y_i \cup Y_{n+i}$  and  $Y_j \cup Y_{n+j}$  are  $K$ -element sets. Furthermore,

$$(Y_i \cup Y_{n+i}) \cap (Y_j \cup Y_{n+j}) = (Y_i \cap Y_j) \cup (Y_i \cap Y_{n+j}) \cup (Y_{n+i} \cap Y_j) \cup (Y_{n+i} \cap Y_{n+j}).$$

Observe that these four sets are pairwise disjoint, and their cardinalities are  $w_{ij} = u_{ij}$ ,  $w_{i,j+n} = v_{ij}$ ,  $w_{i+n,j} = v_{ji}$  and  $w_{i+n,j+n} = K - u_{ij} - v_{ij} - v_{ji}$ , respectively. Thus, the cardinality of the set  $(Y_i \cup Y_{n+i}) \cap (Y_j \cup Y_{n+j})$  is  $K$ , and since the cardinality of each  $Y_i \cup Y_{n+i}$  and  $Y_j \cup Y_{n+j}$  is  $K$ , we conclude  $Y_i \cup Y_{n+i} = Y_j \cup Y_{n+j}$ . This finally shows that  $U$  has cardinality  $K$ , and thus each  $Y_{n+i}$  is  $U \setminus Y_i$  for all  $i \leq n$ . This immediately gives us a set representation for  $\mathbf{U}, \mathbf{V}$ .

To conclude, once we guessed a bounded solution to  $\Psi(D, \Sigma)$  (all components are at most  $M$ ), we proceed to compute in polynomial time the matrix  $\mathbf{W}$  from  $\mathbf{U}$  and  $\mathbf{V}$ , and then run a non-deterministic polynomial time algorithm on it to check if  $\mathbf{W}$  is an intersection pattern. Putting everything together, we see that this nondeterministic polynomial time algorithm returns ‘yes’ iff there is a bounded solution (and thus, there is a solution) to  $\Psi(D, \Sigma)$  for which  $\mathbf{U}, \mathbf{V}$  have a set representation. By Lemma 5.2, this happens if and only if there exists an XML tree  $T$  such that  $T \models D$  and  $T \models \Sigma$ .

This completes the proof of Theorem 5.1. □

We next investigate implication problems.

**THEOREM 5.4.** *For each of  $\mathcal{C}_{K,IC}^{Unary}$  and  $\mathcal{C}_{K,FK}^{Unary}$ , the implication problem is coNP-complete, even under the primary key restriction.*

*Proof:* The implication problem for  $\mathcal{C}_{K,IC}^{Unary}$  is to determine, for a DTD  $D$ , a set  $\Sigma$  of  $\mathcal{C}_{K,IC}^{Unary}$  constraints, and a constraint  $\varphi$  (unary key or unary inclusion constraint), whether  $(D, \Sigma) \vdash \varphi$ . Note that  $(D, \Sigma) \vdash \varphi$  iff there is no XML tree  $T$  with  $T \models D \wedge \bigwedge \Sigma \wedge \neg \varphi$ , and  $\Sigma \cup \{\neg \varphi\}$  is a set of  $\mathcal{C}_{K^-,IC^-}^{Unary}$  constraints. Thus by Theorem 5.1, the implication problem for  $\mathcal{C}_{K,IC}^{Unary}$  is in coNP. One can show that it is coNP-hard under the primary key restriction using an argument similar to the proof of Theorem 4.10. Similarly for the implication problem for  $\mathcal{C}_{K,FK}^{Unary}$ .  $\square$

Finally, along the same lines as Corollary 4.11, we show the following:

**COROLLARY 5.5.** *For a fixed DTD, the following problems can be determined in PTIME:*

- The implication problem for  $\mathcal{C}_{K,FK}^{Unary}$ .
- The consistency problem for  $\mathcal{C}_{K^-,IC^-}^{Unary}$ .

*Proof:* Let  $D$  be a DTD and  $\Sigma$  be a set of  $\mathcal{C}_{K^-,IC^-}^{Unary}$  constraints over  $D$ . Let  $\Psi'(D, \Sigma)$  be the system of linear inequalities determined by  $D$  and  $\Sigma$ , as defined in the proof of Theorem 5.1. As in the proof of Corollary 4.11, one can show that the number of variables in  $\Psi'(D, \Sigma)$  is bounded by a function on the size of  $D$ . Therefore, when  $D$  is fixed, the number of variables in  $\Psi'(D, \Sigma)$  is bounded by a constant. It is known that when the number of variables in a system of linear inequalities is bounded, it can be determined in PTIME whether the system admits an integer solution [29]. By the proofs of Lemma 5.2 and Theorem 5.1,  $\Psi'(D, \Sigma)$  admits an integer solution if and only if there is an XML tree  $T$  such that  $T \models D$  and  $T \models \Sigma$ . Thus Corollary 5.5 follows from Theorems 5.1 and 5.4.  $\square$

## 6. CONCLUSION

We have studied the consistency problems associated with four classes of integrity constraints for XML. We have shown that in contrast to its trivial counterpart in relational databases, the consistency problem is undecidable for  $\mathcal{C}_{K,FK}$ , the class of multi-attribute keys and foreign keys. This demonstrates that the interaction between DTDs and key/foreign key constraints is rather intricate. This negative result motivated us to study  $\mathcal{C}_{K,FK}^{Unary}$ , the class of unary keys and foreign keys, which are commonly used in practice. We have developed a characterization of DTDs and unary constraints in terms of linear integer constraints. This establishes a connection between DTDs, unary constraints and linear integer programming, and allows us to use techniques from combinatorial optimization in the study of XML constraints. We have shown that the consistency problem for  $\mathcal{C}_{K,FK}^{Unary}$  is NP-complete. Furthermore, the problem remains in NP for  $\mathcal{C}_{K^-,IC^-}^{Unary}$ , the class of unary keys, unary inclusion constraints and their negations.

We have also investigated the implication problems for XML keys and foreign keys. In par-

	<i>multi-attribute keys, foreign keys</i>	<i>unary keys, foreign keys</i>	<i>primary, unary keys, foreign keys</i>	<i>DTD fixed, unary keys, foreign keys</i>	<i>multi-attribute keys only</i>
<i>consistency</i>	undecidable	NP-complete	NP-complete	PTIME	linear time
<i>implication</i>	undecidable	coNP-complete	coNP-complete	PTIME	linear time

Fig. 5. The main results of the paper

ticular, we have shown that the problem is undecidable for  $\mathcal{C}_{K,FK}$  and it is coNP-complete for  $\mathcal{C}_{K,FK}^{Unary}$  constraints. Several PTIME decidable cases of the implication and consistency problems have also been identified. The main results of the paper are summarized in Figure 5.

It is worth remarking that the undecidability and NP-hardness results also hold for other schema specifications beyond DTDs, such as XML Data [27], XML Schema [36] and the generalization of DTDs proposed in [33]. It remains open, however, whether the upper bounds (i.e., the decidability and NP membership results) are still intact in these settings.

This work is a first step towards understanding the interaction between DTDs and integrity constraints. A number of questions remain open. First, we have only considered keys and foreign keys defined with XML attributes. We expect to extend techniques developed here for more general schema and constraint specifications. Second, other constraints commonly found in databases, e.g., inverse constraints, deserve further investigation. Third, a lot of work remains to be done on identifying tractable yet practical classes of constraints and on developing heuristics for consistency analysis. Finally, a related project is to use integrity constraints to distinguish good XML design (specification) from bad design, along the lines as normalization of relational schemas. Coding with linear integer constraints gives us decidability for some implication problems for XML constraints, which is a first step towards a design theory for XML specifications.

**Acknowledgments.** We thank Michael Benedikt, Alberto Mendelzon, Frank Neven and Dan Suciu for helpful discussions. We are grateful to the referees for valuable suggestions on simplifying the proofs and on improving the paper.

## REFERENCES

- [1] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- [2] S. Abiteboul and V. Vianu. Regular path queries with constraints. *JCSS*, 58(4):428–452, 1999.
- [3] V. Apparao et al. Document Object Model (DOM) Level 1 Specification. W3C Recommendation, Oct. 1998. <http://www.w3.org/TR/REC-DOM-Level-1/>.

- [4] C. Baru, A. Gupta, B. Ludäscher, R. Marciano, Y. Papakonstantinou, P. Velikhov, and V. Chu. XML-based information mediation with MIX. In *SIGMOD'99*, pages 597–599.
- [5] C. Beeri and T. Milo. Schemas for integration and translation of structured and semi-structured data. In *ICDT'99*, pages 296–313.
- [6] T. Bray, J. Paoli, and C. M. Sperberg-McQueen. Extensible Markup Language (XML) 1.0. W3C Recommendation, Feb. 1998. <http://www.w3.org/TR/REC-xml/>.
- [7] P. Buneman, S. Davidson, W. Fan, C. Hara, and W. Tan. Keys for XML. In *WWW'10*, 2001.
- [8] P. Buneman, S. Davidson, W. Fan, C. Hara, and W. Tan. Reasoning about keys for XML. In *DBPL'01*, 2001.
- [9] P. Buneman, W. Fan, and S. Weinstein. Interaction between path and type constraints. In *PODS'99*, pages 56–67.
- [10] P. Buneman, W. Fan, and S. Weinstein. Path constraints in semistructured databases. *JCSS*, 61(2):146–193.
- [11] D. Calvanese, G. De Giacomo, and M. Lenzerini. Representing and reasoning on XML documents: A description logic approach. *J. Logic and Computation*, 9(3):295–318, 1999.
- [12] D. Calvanese and M. Lenzerini. Making object-oriented schemas more expressive. In *PODS'94*, pages 243–254.
- [13] D. Calvanese and M. Lenzerini. On the interaction between ISA and cardinality constraints. In *ICDE'94*, pages 204–213.
- [14] D. Chamberlin et al. XQuery 1.0: An XML Query Language. W3C Working Draft, June 2001. <http://www.w3.org/TR/xquery>.
- [15] J. Clark. XSL Transformations (XSLT). W3C Recommendation, Nov. 1999. <http://www.w3.org/TR/xslt>.
- [16] J. Clark and S. DeRose. XML Path Language (XPath). W3C Recommendation, Nov. 1999. <http://www.w3.org/TR/xpath>.
- [17] S. S. Cosmadakis, P. C. Kanellakis, and M. Y. Vardi. Polynomial-time implication problems for unary inclusion dependencies. *J. ACM*, 37(1):15–46, Jan. 1990.
- [18] W. Fan and L. Libkin. On XML integrity constraints in the presence of DTDs. In *PODS'01*, pages 114–125.
- [19] W. Fan and J. Siméon. Integrity constraints for XML. In *PODS'00*, pages 23–34.
- [20] M. Fernandez, D. Florescu, A. Levy, and D. Suciu. Verifying integrity constraints on Web sites. In *IJCAI'99*, pages 614–619.
- [21] D. Florescu, L. Raschid, and P. Valduriez. A methodology for query reformulation in CIS using semantic knowledge. *Int'l J. Cooperative Information Systems (IJCIS)*, 5(4):431–468, 1996.
- [22] M. Garey and D. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman and Company, 1979.
- [23] C. Hara and S. Davidson. Reasoning about nested functional dependencies. In *PODS'99*, pages 91–100.
- [24] J. E. Hopcroft, R. Motwani, and J. D. Ullman. *Introduction to Automata Theory, Languages and Computation* (2nd Edition). Addison Wesley, 2000.
- [25] M. Ito and G. E. Weddell. Implication problems for functional constraints on databases supporting complex objects. *JCSS*, 50(1):165–187, 1995.
- [26] P. C. Kanellakis. On the computational complexity of cardinality constraints in relational databases. *Information Processing Letters*, 11(2):98–101, Oct. 1980.
- [27] A. Layman et al. XML-Data. W3C Note, Jan. 1998. <http://www.w3.org/TR/1998/NOTE-XML-data>.

- [28] D. Lee and W. W. Chu. Constraint-preserving transformation from XML document type to relational schema. In *ER'00*, pages 323–338, 2000.
- [29] H. W. Lenstra. Integer programming in a fixed number of variables. *Math. Oper. Res.*, 8:538–548, 1983.
- [30] J. Melton and A. Simon. *Understanding the New SQL: A Complete Guide*. Morgan Kaufman, 1993.
- [31] F. Neven. Extensions of attribute grammars for structured document queries. In *DBPL'99*, pages 99–116, 1999.
- [32] C. H. Papadimitriou. On the complexity of integer programming. *J. ACM*, 28(4):765–768, 1981.
- [33] Y. Papakonstantinou and V. Vianu. Type inference for views of semistructured data. In *PODS'00*, pages 35–46.
- [34] L. Popa. *Object/Relational Query Optimization with Chase and Backchase*. PhD thesis, University of Pennsylvania, 2000.
- [35] J. Robie, J. Lapp, and D. Schach. XML Query Language (XQL). Workshop on XML Query Languages, Dec. 1998.
- [36] H. S. Thompson et al. XML Schema. W3C Working Draft, May 2001. <http://www.w3.org/XML/Schema>.
- [37] J. D. Ullman. *Database and Knowledge Base Systems*. Computer Science Press, 1988.
- [38] J. Widom. Data management for XML: Research directions. In *IEEE Data Eng. Bull.*, 22(3): 44-52, 1999.