

# Querying Regular Sets of XML Documents

Slawomir Staworko, Emmanuel Filiot, Jan Chomicki

# ▶ To cite this version:

Slawomir Staworko, Emmanuel Filiot, Jan Chomicki. Querying Regular Sets of XML Documents. International Workshop on Logic in Databases (LiD), May 2008, Rome, Italy. 2008. <inria-00275491>

# HAL Id: inria-00275491 https://hal.inria.fr/inria-00275491

Submitted on 24 Apr 2008  $\,$ 

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Querying Regular Sets of XML Documents<sup>\*</sup>

Slawomir Staworko<sup>1,2</sup>, Emmanuel Filiot<sup>1</sup>, and Jan Chomicki<sup>2</sup>

 <sup>1</sup> INRIA Lille Nord-Europe
<sup>2</sup> Department of Computer Science and Engineering University at Buffalo, Buffalo, NY

Abstract. We investigate the problem of querying (regular) sets of XML documents represented with tree automata and we consider n-ary tree automata queries whose expressive power captures MSO on trees. Because finite automata can represent infinite sets of documents, we propose the notions of universal and existential query answers, answers that are present resp. in all and some documents. We study complexity of query answering and show that computing existential query answers is in PTIME if we assume the arity of the query to be a fixed parameter. On the other hand, computing universal query answers is EXPTIMEcomplete, but we show that it is in PTIME if we assume that the query is fixed (data complexity). Finally, we argue that the framework captures problems central to many novel XML applications like querying inconsistent XML documents. In particular, we demonstrate how to use our framework to compute consistent query answers in XML documents that do not satisfy the schema. This solution significantly extends our previous results in this area.

# 1 Introduction

In this paper we investigate the problem of querying potentially large sets of XML documents having a small compact representation by a finite tree automaton. Our work is inspired by the problem of evaluating a query in a document that possibly does not satisfy the schema. Since the satisfaction of the schema is usually assumed during formulation of a query, evaluation of the query against a document that does not satisfy the schema may yield incorrect and misleading answers. This problem has been previously recognized in the setting of relational databases [3] and the proposed framework of repairs and consistent query answers (CQA) has been adapted to semi-structured databases [18, 10, 11]. A *repair* is a document satisfying the schema and obtained from the original document by a minimal number of standard edit operations [1, 5]: inserting, deleting, and modifying an element of the document. An answer is *consistent* (also called *valid*) if it is an answer to the query in every repair.

Our research shows that the set of repairs of a document is a regular language that has a compact representation by a finite (weighted) tree automata.

<sup>\*</sup> Research supported by NSF Grants IIS-0119186 and IIS-0307434.

This reduces the problem of querying the set of repairs to a more general problem of querying a regular set of documents represented by a tree automaton. Because the set of documents represented by an automaton can be very large, an approach where we return the collection of the sets of answers in every tree may be simply inappropriate for many applications. Consequently, we propose two ways of computing answers to queries in sets of documents: (i) universal answers that are present in every document and (ii) existential answers that are present in some document. Obviously, universal answers capture consistent answers. We also note that the existential answers capture the notion of *possible* answers, i.e. the answers that are present in some repair [10], often considered next to consistent answers in the framework of CQA. Another motivation to study the problem of querying regular sets of documents is that there are other frameworks where the queries are evaluated not on the instance itself but rather on the set of instances obtained by some (often nondeterministic) process from the original one. An example is XML data exchange [4]: the queries are formulated against *target* schema whose instances (called *solutions*) are obtained from a source instance with a set of source-to-target dependencies specifying how the parts of the source instance translate to an instance in target schema. In the data exchange scenario, the notion of universal and existential answers coincide resp. with *certain* and *maybe* answers [15]. It is, however, yet to be seen if the set of solutions can be represented by a tree automaton.

Now, we briefly summarize our contributions:

- We define the problem of universal and existential querying of sets of documents represented by a finite tree automaton with attribute values and we consider *n*-ary queries expressed with tree automata.
- We thoroughly study the complexity of computing existential and universal answer. For computing existential answers, we show that its combined complexity is NP-complete, but its complexity parametrized by the arity of the query is FPT (*Fixed Parameter Tractable*) [9]. This result is not surprising as the number of answers to an *n*-ary query can be exponential in *n*. Computing universal answers is, however, EXPTIME-complete in terms of both the combined and parametric complexity. On the other hand we show that its data-complexity is PTIME.
- We show how to compute consistent query answers by constructing a *repair* automaton that defines the set of all repairs of a document w.r.t. a schema (expressed with a tree automaton). This extends our previous results [18, 17] in several directions: (i) we consider *n*-ary automata queries whose expressive power captures MSO as compared to a restricted class of unary Core XPath queries, (ii) schema can be expressed using tree automata which are strictly more expressive than DTDs, (iii) we show that data complexity of computing consistent answer to any *n*-ary automata query is PTIME, (iv) we consider a more general set of editing tree operations that allow to operate on inner nodes as compared to operations on leaves only.

**Related work** [10] investigates querying XML documents that are valid but violate functional dependencies. Two repairing actions are considered: updating element values with a *null* value and marking nodes as unreliable. Such nodes are simply omitted in the query answers. Only simple descending path queries  $a_1/a_2/\ldots/a_n$  are considered. A polynomial algorithm for computing consistent and possible answers is presented. [11] considers editing operations operating on leaves only to define the set of repairs for consistent querying of documents that violate functional dependencies. A different, set-theoretic, notion of minimality is used when defining repairs. Both consistent and possible answers are considered. For restricted classes of functional dependencies and DTDs a polynomial algorithm is proposed to compute consistent answers to *n*-ary conjunctions of path expressions.

[12] considers evaluation of monadic Datalog queries on compressed trees (represented by a tree automaton) and shows that combined complexity is PSPACE-complete and data complexity is PTIME. We note that this framework is close to existential querying: a (finite) set of trees can be gathered in one tree with a new root symbol. The difference between the complexity results (for unary queries we have polynomial algorithm in terms of combined and data complexity) comes from different representations of queries. [16] extends this approach to trees represented by straight-line context-free grammars which is strictly stronger than regular languages.

[21] study the problem of checking if a document is within a specified *alignment distance* to the given schema. We note that the edit distance is more general than the alignment distance which imposes certain conditions on the sequence of editing operations [6] and hence our approach is more general. A compact representation of all repairs (obtained with restricted sequences of editing operations) as a regular language is also presented.

The paper is organized as follows. Section 2 contains basic XML notions and streaming tree automata. In Section 3 we define existential and universal answers to n-ary queries and present algorithm for computing them. In Section 4 we study computational implications of our framework. Section 5 shows how to use our framework to compute consistent query answers. Because of space limitations the proofs are omitted; they can be found in the appendix available at [19].

## 2 Basic notions

### 2.1 Trees and streams

We model XML documents using ordered unranked trees whose nodes are labeled with elements of a finite set of symbols  $\Sigma$ . Every node is additionally labeled with an attribute whose value is drawn from an infinite set  $\Lambda$ . We denote the set of all trees by  $\mathcal{T}$ . The size of t, denoted by |t|, is the number of nodes of t.

In our framework only the attribute values are used to define query answers. The attributes of a tree can store unique node identifiers and we call such trees *standard.* In general, however, the attribute can store (possibly repeated) data values.

In this paper we work mainly with the serialized version of trees, i.e. wellformed sequences of opening and closing tags (corresponding to a preorder traversal of the tree) with attribute values associated to the opening tags. The set of all tags is  $\Sigma_{\Diamond} = \{\langle a \rangle | a \in \Sigma\} \cup \{\langle /a \rangle | a \in \Sigma\}$ . When working with a serialized version of a tree  $e_1, \ldots, e_n$  we write  $tag(e_i)$  for the tag of  $e_i$  and  $att(e_i)$  for the attribute value of  $e_i$ . Given a tree t its serialized version is denoted by  $\bar{t}$ . We



**Fig. 1.** An example of a tree  $t_0$  and the corresponding tag sequence

use also unranked terms over the signature  $\Sigma \times \Lambda$  to represent trees. Figure 1 contains an example of a tree  $t_0 = c^{l_0}(a^{l_1}, b^{l_2}, a^{l_3}, b^{l_4}, a^{l_5})$  and its serialization.

### 2.2 Streaming tree automata

To capture regular tree languages we use streaming tree automata [13] which are Visibly Pushdown Automata [2] working on serializations of unranked trees. They are allow to capture Extended DTDs [14] which makes them equivalent to standard (ranked) tree automata working on encodings of unranked trees [8]. We choose this model because it is better fitted to capture repairs of XML documents: repairs are obtained by edit operations on nodes and those easily translate to string edit operations on pairs of matching tags. We also extend this model by allowing it to specify attribute values: this way the repairs can be defined in terms of the attribute values of the original document. For simplicity of presentation we fix the set of attribute values  $\Lambda$ .

**Definition 1.** An attributed streaming tree automaton (attributed STA) is a tuple  $M = (\Sigma, \Gamma_M, Q_M, I_M, \Delta_M, F_M)$ , where  $\Sigma$  is a finite set of node labels,  $\Gamma_M$  is a finite set of stack symbols,  $Q_M$  is a finite set of states, and  $I_M \subseteq Q_M$  is the set of initial states,  $F_M \subseteq Q_M$  is the set of final states.  $\Delta_M$  is a finite set of transitions of one of the following types: opening transition  $p \xrightarrow{\langle a \rangle: \gamma} q$ , and closing transition  $p \xrightarrow{\langle A \rangle: \gamma} q$ , where  $p, q \in Q_M$ ,  $a \in \Sigma$ ,  $\gamma \in \Gamma_M$ , and  $l \in \Lambda \cup \{*\}$  (\* is a wildcard).

The size of M is  $|M| = |Q_M| + |\Delta_M|$  and by  $Dom(M) \subseteq \Lambda \cup \{*\}$  we denote the set of different attribute labels used in  $\Delta_M$ .

Essentially, an attributed STA is a push-down automaton working on sequences of tags with the stack manipulation restricted to: placing a symbol on the stack when reading an opening tag, and removing top symbol from the stack when reading a closing tag. A *configuration* is an tuple  $(q, \bar{\alpha}, \bar{s})$ , where q is the current state,  $\bar{\alpha}$  is the current stack, and  $\bar{s}$  is the remaining (possibly unbalanced) tag sequence. The *move* relation  $\rightarrow_M$  is a binary relation on configurations defined as follows:

- (i)  $(q, \bar{\alpha}, e \cdot \bar{s}) \to_M (p, \gamma \cdot \bar{\alpha}, \bar{s})$  if  $q \xrightarrow{\langle a \rangle : \gamma}{l} p \in \Delta_M$ ,  $tag(e) = \langle a \rangle$ , and if  $l \neq *$ , then att(e) = l.
- (ii)  $(q, \gamma \cdot \bar{\alpha}, e \cdot \bar{s}) \to_M (p, \bar{\alpha}, \bar{s})$  if  $q \xrightarrow{\langle /a \rangle : \gamma} p \in \Delta_M$  and  $tag(e) = \langle /a \rangle$ .

 $\rightarrow_M^*$  is the reflexive and transitive closure of  $\rightarrow_M$ . The tree language of M is defined as  $L(M) = \{t \mid (q, \varepsilon, \overline{t}) \rightarrow_M^* (p, \varepsilon, \varepsilon), q \in I_M, p \in F_M\}.$ 

An *STA* is an attributed STA that imposes no restrictions on the attribute values, i.e. it uses only \*. Then, we also omit the attribute labels altogether. Figure 2 contains an example of an STA  $M_0$  that recognizes trees with root label c satisfying the DTD  $D_0$  given by the rules:  $c \to (a \cdot b)^* \cdot a, a \to \epsilon, b \to \epsilon$ . We observe that the tree  $t_0$  (Fig. 1) satisfies  $M_0$ .



**Fig. 2.** The STA  $M_0$  for the DTD  $D_0$ 

Weighted STAs To represent the sets of minimal repairs we further extend attributed STAs by assigning to every transition its *weight*, a non-negative real value. The weights are used to restrict the set of recognized trees to those with a run of minimal summary weight.

More formally, a weight of a run is the sum the weights of the transition used at each step. For a tree (with an accepting run) we associate the minimal weight of its accepting run. L(M) contains only the trees that whose weight is equal to the minimum of weights of all trees (with an accepting run).

### 2.3 STA queries

We define an *n*-ary query as a function that takes a tree and returns a set of n-ary tuples of values from  $\Lambda$  (that are used in the tree).

To define an *n*-ary query we use an extension of (standard) STAs where with every opening transition we associate a set of variables  $X \subseteq \{x_1, \ldots, x_n\}$  that indicates the positions of the resulting tuple  $(x_1,\ldots,x_n)$  that are to be filled when the transition is used in a run (the positions are filled with the attribute value). Each position of the resulting tuple has to be filled *exactly* once during a run, otherwise the run is non-accepting. The set of query answers consists of all tuples obtained from all accepting runs. We use name STA queries to refer to such automata and to distinguish them from STAs we use Greek capital letters  $\Phi, \Psi, \ldots$  We also put the sets X in the superscript of the opening transitions. Figure 3 contains an STA query  $\Phi$  selecting pairs  $(x_1, x_2)$  of any node a and its immediate right sibling b (satisfaction of DTD  $D_0$  is assumed). On the tree  $t_0$ 

Fig. 3. An example of a binary STA query

(Fig. 1) this query has two answers:  $(l_1, l_2)$  and  $(l_3, l_4)$ .

We define query answers formally as follows. A configuration of an n-ary STA query  $\Phi$  is an element  $(q, \bar{\alpha}, \bar{s}, \tau)$ , where  $q, \bar{\alpha}$ , and  $\bar{t}$  are as before, and  $\tau \in (\Lambda \cup \{\bot\})^n$  is the tuple of values assigned so far, with  $\bot$  (the *null* value) indicating that the value has not been yet assigned. We define the move relation analogously:

- (i)  $(q, \bar{\alpha}, e \cdot \bar{s}, \tau) \to_{\varPhi} (p, \gamma \cdot \bar{\alpha}, \bar{s}, \tau')$  if  $q \xrightarrow{\langle a \rangle : \gamma} p \in \Delta_{\varPhi}, tag(e) = \langle a \rangle, \tau_i = \bot$  for every  $x_i \in X$ , and  $\tau' = \tau[X/att(e)]$ , (ii)  $(q, \gamma \cdot \bar{\alpha}, e \cdot \bar{s}, \tau) \to_{\varPhi} (p, \bar{\alpha}, \bar{s}, \tau)$  if  $q \xrightarrow{\langle a \rangle : \gamma} p \in \Delta_{\varPhi}$  and  $tag(e) = \langle a \rangle$ .

Again,  $\rightarrow_{\Phi}^*$  is the reflexive and transitive closure of  $\rightarrow_{\Phi}$ . The set of answers to an *n*-ary STA query  $\Phi$  in t is  $QA(\Phi, t) = \{\tau \in \Lambda^n | (q, \varepsilon, \overline{t}, (\bot, ..., \bot)) \rightarrow^*_{\Phi}$  $(p,\varepsilon,\varepsilon,\tau), p \in I_{\Phi}, q \in F_{\Phi} \}.$ 

It is known [13] that over standard trees n-ary STA queries have the same expressive power as MSO formulas with n free variables over the first-child, next-sibling signature of unranked trees In particular, STA queries subsume unary CoreXPath queries. It should be noted, however, that similarly translating an MSO formulas may yield a automata of non-elementary size [20, 13]. This generally applies also to CoreXPath queries, but it can be easily seen that simple descending XPath queries with no test expressions translate to STA queries of linear size.

#### 3 Existential and universal querying of attributed STAs

Now, we consider querying sets of trees defined by attributed STAs. First, we note that the set of trees defined by an attributed STA may be infinite and even if the STA is weighted, the number of trees may be exponential in the size of the automaton (modulo different attribute values). Therefore, an approach where we return the collection of the sets of answers obtained in every tree may be inappropriate for many applications. Consequently, we propose two ways of querying sets of trees.

**Definition 2 (Universal and existential answers).** Given a (possibly weighted) attributed STA M such that  $L(M) \neq \emptyset$  and n-ary STA query  $\Phi$ 

- the universal answers to  $\Phi$  in M are  $QA^{\forall}(\Phi, M) = \bigcap_{t \in L(M)} QA(\Phi, t)$ ,
- the existential answers to  $\Phi$  in M are  $QA^{\exists}(\Phi, M) = \bigcup_{t \in L(M)} QA(\Phi, t)$ .

We note that the language defined by a weighted automaton is empty if and only if the corresponding automaton without weights defines an empty language. Emptiness of an attributed STA can be tested in cubic time using the classical algorithm for PDAs. Because our algorithms have complexity of a higher degree, from now on we will assume that we always deal with automata defining nonempty language. Also, the set  $QA^{\forall}(\Phi, M)$  is always finite, but  $QA^{\exists}(\Phi, M)$ may be infinite if M uses wildcards. Hence, we also allow to use wildcards in answers to finitely represent  $QA^{\exists}$ .

Now, we present Algorithm 1 computing existential answers to an *n*-ary query  $\Phi$  in an attributed STA *M*. This algorithm is based on a product technique

<b>Algorithm 1</b> Computing existential answers to <i>n</i> -ary $\Phi$ in $M$
function $QA^{\exists}(\Phi, M)$
<b>macros</b> : $Q :\equiv Q_{\Phi} \times Q_M, I :\equiv I_{\Phi} \times I_M, F :\equiv F_{\Phi} \times F_M$
$(p_1,q_1) \xrightarrow{x:(\gamma_1,\gamma_2)}_{X:=l} (p_2,q_2) :\equiv p_1 \xrightarrow{x:\gamma_1}_X p_2 \in \Delta_{\varPhi} \land q_1 \xrightarrow{x:\gamma_2}_l q_2 \in \Delta_M,$
1: for $(u,v) \in Q^2$ do
2: $T_0[u, v] = \begin{cases} \{(\bot, \dots, \bot)\}, & \text{if } u = v, \end{cases}$
$ [\varnothing, otherwise.] $
3: for $i \leftarrow 1, \dots, n Q ^2$ do
4: for $(u,v) \in Q^2$ do
5: $H_i[u, v] = T_{i-1}[u, v]$
6: for $j \leftarrow 1, \dots, \lceil \log(n Q ^2) \rceil$ do
7: for $(u,v) \in Q^2$ do
8: $H_i[u,v] \leftarrow H_i[u,v] \cup \bigcup \{merge(H_i[u,w],H_i[w,v])   w \in Q\}$
9: for $(u,v) \in Q^2$ do
10: $T_{i}[u,v] \leftarrow T_{i-1}[u,v] \cup \bigcup \{assign_{X}(H_{i}[u',v'],l)   u \xrightarrow{\langle a \rangle : \gamma} X := l} u' \wedge v' \xrightarrow{\langle /a \rangle : \gamma} v \}$
11: <b>return</b> $\{\tau \in (\Lambda \cup \{*\})^n   \tau \in T_{n Q ^2}[u, v] \land u \in I \land v \in F\}$
end function

of the two input automata. Essentially, it evaluates the query on every tree of height and width  $\leq n|Q|^2$ , where  $Q = Q_{\Phi} \times Q_M$ . This procedure yields correct results thanks to pumping properties of STAs. In particular, if there is a tree

 $t \in L(M)$  and tuple  $\tau \in QA^{\exists}(\Phi, t)$ , there is also a tree  $t^*$  whose depth and width is bounded by  $n|Q|^2$ , such that  $t^* \in L(M)$  and  $\tau \in QA^{\exists}(\Phi, t^*)$ . Consequently, we need to consider query runs of depth and width bound by  $n|Q|^2$ . This space can be explored with a simple dynamic programming technique because runs of an STA on a tree share the structure of the tree. In particular,  $T_i$  and  $H_i$  store all tuples "collected" from runs on resp. trees and hedges (sequences of trees) of depth  $\leq i$  and width  $\leq n|Q|^2$ . We use  $assign_X(A, l)$  to assign the value lon positions X to every tuple from A (tuples having a value different from  $\perp$ on those positions are discarded). merge(A, B) returns the set of merged tuples from sets A and B (two tuples having assigned value on the same position cannot be merged). An easy complexity analysis shows that:

**Theorem 1.** For an attributed STA M and n-ary query  $\Phi$  Algorithm 1 computes  $QA^{\exists}(\Phi, M)$  in time  $O((|\Phi||M|)^{6}|Dom(M)|^{2n})$ .

Extending the algorithm to weighted attributed STAs is not difficult because minimal runs enjoy optimal substructure properties. Also, if we first perform the run of the algorithm on the attributed STA where every attribute value has been replaced by one unique constant, we can find which tuples in the intermediate steps are removed by the *apply* and *merge* operations. This allow us to replace the  $|Dom(M)|^{2n}$  factor by  $2^{2n}|QA^{\exists}(\Phi, M)|$ .

**Corollary 1.** For any weighted attributed STA M and n-ary query  $\Phi$  the set  $QA^{\exists}(\Phi, M)$  can be computed in time  $O(2^{2n}(|\Phi||M|)^{6}|QA^{\exists}(\Phi, M)|)$ .

We observe that the high complexity in terms of |M| comes from the particular pumping properties of STAs. We note, however, that the algorithm could be easily adapted to standard tree automata working on binary representation of unranked trees. Those automata enjoy nicer pumping properties and in particular the complexity would be cubic in terms of the size of the input automata.

To compute universal answers we observe that  $QA^{\forall}(\Phi, M) \subseteq QA^{\exists}(\Phi, M)$ . Instead of computing universal answers directly, we compute  $QA^{\exists}(\Phi, M)$  and on every tuple we perform *tuple check*, i.e. we find if the tuple is an answer in every tree. Because we use it as a tool for analyzing the complexity of universal and existential query answers, we formally define it:

**Existential (universal) tuple check** is a decision problem where given an attributed STA M, an *n*-ary STA query  $\Phi$ , and a tuple  $\tau \in (\Lambda \cup \{*\})^n$ find if  $\tau$  is an existential (universal resp.) answer to  $\Phi$  in M.

**Theorem 2.** Universal tuple check can be decided in time  $O(f(|\Phi|)|M|^3))$ , where  $f(|\Phi|) = 2^{|\Phi|^2 2^{2n}}$  (and  $O(f(|\Phi|) \log(|M|)|M|^4)$ ) if M is weighted).

Consequently, we obtain a characterization of the simple procedure of computing universal answers.

**Corollary 2.** For any weighted attributed STA M and n-ary query  $\Phi$  the set  $QA^{\forall}(\Phi, M)$  can be computed in time  $O(2^{|\Phi|^2 2^{2n}} |\Phi|^6 |M|^6 |QA^{\exists}(\Phi, M)|)$ .

# 4 Complexity analysis

Now, we analyze tractability of computing existential and universal query answers by investigating the complexity of universal and existential tuple check (defined in the previous section). We start with combined complexity where all the elements are considered to be the part of the input.

**Theorem 3.** Combined complexity of existential and universal tuple check are NP-complete and EXPTIME-complete respectively.

We remark that the EXPTIME-hardness is proved with a reduction of the containment problem of two tree automata to a universal tuple check where one of the automata is treated as a 0-ary query.

Next, we observe that if the arity of the query is fixed, then the existential tuple check can be done in polynomial time. Moreover, the degree of the polynomial does not depend on the arity of the query. Hence, we can characterize the (multiplicative) *fixed parametric complexity* [9].

**Corollary 3.** When the arity of the STA query is a parameter, the existential tuple check is FPT (Fixed Parameter Tractable).

We note that universal tuple check remains intractable when fixing the arity of the query as the EXPTIME-hardness proof uses an STA query of arity 0.

Finally, Theorem 2 give us a characterization of data complexity [22] of universal query answers (the query is assumed to be fixed).

Corollary 4. Data complexity of universal tuple check is PTIME.

# 5 Consistent querying of XML documents

We recall the basic notions of the framework of consistent query answers for semistructured databases. The process of repairing an XML document is modeled with the standard edit operations on trees: (i) renaming the node, (ii) deleting a node (different than the root) which involves promoting its children to the parent of the node (placed in the same order from the position of the node), and (iii) inserting a node (different than the root) with a possible adoption of a list of subsequent children from the parent of the node (dual to the delete operation). With every editing operation we associate a cost:  $c_R$ ,  $c_D$ , and  $c_I$  the costs for renaming, deleting, and inserting a node respectively. We note, however, that our approach can be easily extended to weights that depend on properties of the node, for example its label.

The *edit distance*  $d(t_1, t_2)$  between two trees  $t_1$  and  $t_2$  is the minimal cost of transforming  $t_1$  to  $t_2$  with a series of edit operations. Given a tree t and an STA M (expressing the schema), we define the distance between t and M, denoted d(M, t), as the minimum edit distance between t and any t' valid w.r.t. M.

A repair of t w.r.t. M is a tree  $t' \in L(M)$  such that d(t,t') = d(t,M). By Rep(t,M) we denote the set of all repairs of t w.r.t. M. Given an n-ary STA query  $\Phi$  we say that a tuple  $\tau$  is a consistent (or valid) answer to  $\Phi$  in t w.r.t. M if and only if  $\tau$  is an answer to  $\Phi$  in every repair of t w.r.t. M. By  $CQA(\Phi,t,M)$  we denote the set of all consistent answers to  $\Phi$  in t w.r.t. M.

#### 5.1**Repair automaton**

In this part we define a weighted attributed STA that defines the set of all repairs of a document. For ease of construction we allow the use of  $\epsilon$ -transitions. As they have no attribute value and perform no operations on the stack, we can easily remove them by standard closure (remembering to aggregate the weight). Also, we make a natural assumption that the schema does not allow an empty tree.

**Definition 3.** Let M be an STA, t be a tree with n nodes, and  $\overline{t} = (e_1, \ldots, e_{2n})$ be the serialization of t. Assume that the nodes of t are numbered with consecutive natural numbers  $0, 1, \ldots, n-1$  in the standard document order. For an opening or closing tag  $e_i$  let  $m_i$  be the number assigned to that node.

The repair automaton of t w.r.t. M is a weighted attributed STA R(t, M) = $(\Sigma, \Gamma_R, Q_R, I_R, \Delta_R, F_R)$  where:  $\Gamma_R = \Gamma_M \cup \Gamma_M \times \{1, \dots, n\}, Q_R = I_M \times \{0\} \cup$  $Q_M \times \{1, \ldots, 2n-1\} \cup F_M \times \{2n\}, I_R = I_M \times \{0\}, and F_R = F_M \times \{2n\}.$  The state (q,i) will be denoted as  $q^i$ . The transitions of  $\Delta_R$  (with their attributes and weights) capture edit operations performed on the tag stream as follows:

- $-q^{i-1} \xrightarrow{x:(\gamma,m_i)} p^i$  renaming  $e_i$  to x if  $tag(e_i) \neq x$  and doing nothing if  $att(e_i) \neq x$ 
  - $tag(e_i) = x$ , for every  $i \in \{1, \ldots, 2n\}$  and every  $q \xrightarrow{x:\gamma} p \in \Delta_M$ ; its weight
- is  $c_R/2$  if  $tag(e_i) \neq x$  and 0 if  $tag(e_i) = x$ .  $-q^i \xrightarrow{x:\gamma} p^i$  inserting x (before  $e_i$ ), for every  $i \in \{2, \ldots, 2n-1\}$  and every  $q \xrightarrow{x:\gamma} p \in \Delta_M$ ; its weight is  $c_I/2$  and it is attributed with \* if x is an opening tag. -  $q^{i-1} \stackrel{\varepsilon}{\to} q^i$  deleting  $e_i$ , for every  $i \in \{2, \ldots, 2n-1\}$  and every  $q \in Q_M$ ; its
- weight is  $c_D/2$ .

Example 1. Figure 4 contains the repair automaton of the tree  $t_1 = c^{l_0}(a^{l_1}, b^{l_2})$ (Fig. 1) and the STA  $M_0$  (Fig. 2). The weights are assumed to be  $w_R = w_I =$  $w_D = 1$ . Because this graph is very intricate, for clarity we present in all details only the transitions that produce the minimal trees. Also,  $\gamma_i^j$  is short for  $(\gamma_i, j)$ .  $R(t_1, M_0)$  defines the set of trees  $\{c^{l_0}(a^{l_1})\} \cup \{c^{l_0}(a^{l_1}, b^{l_2}, a^{l_1}) | l \in A\}$ , i.e. the set of repairs of  $t_1$  w.r.t.  $M_0$ .

**Theorem 4.** For any tree t, any STA M, and any STA query  $\Phi$  we have that Rep(t, M) = L(R(t, M)) and  $CQA(\Phi, t, M) = QA^{\forall}(\Phi, R(t, M)).$ 

#### 5.2**Complexity** analysis

From Corollary 2 we get directly:

**Corollary 5.** The data complexity of computing consistent answers to an n-ary query w.r.t. an STA is PTIME.

To further analyse the tractability of consistent query answers we investigate the complexity of the problem of tuple check for consistent query answers [17]. Similarly to universal answers the problem is intractable.

**Theorem 5.** The combined complexity of consistent query answers is  $\Pi_p^2$ -complete if  $w_I > 0$  and EXPTIME-complete if  $w_I = 0$ .



**Fig. 4.** A repair graph  $R(t_1, M_0)$ 

# 6 Conclusions and future work

In this paper we presented the framework of querying regular sets of XML documents. We considered the class of *n*-ary automata queries and introduced notions of *universal* and *existential* answers, i.e. answers present in every and some documents represented by a tree automata. We investigated computational properties of our framework and presented algorithms for computing universal and existential answers. Finally, we used our framework to compute consistent query answers in XML documents that do not satisfy the schema. This solution significantly extends our previous results in this area [18].

We envision several possible directions of further study. Firstly, we would like to investigate the problem of querying regular sets of trees with a query defined in a logical formalism which allow comparison of data-values. Such formalisms have been studied in [7] where FO logic over data trees is considered. It is shown that FO is decidable for the two-variable fragment with a successor relation and a predicate to compare data-values. Another important direction of future study is to investigate if the setting of data exchange could be effectively captured by our framework. Finally, we would like to investigate using more expressive formalisms, for instance context-free tree grammars, to represent the sets of trees.

## References

- S. Abiteboul, J. McHugh, M. Rys, V. Vassalos, and J. L. Wiener. Incremental Maintenance for Materialized Views over Semistructured Data. In *International Conference on Very Large Data Bases (VLDB)*, pages 38–49, 1998.
- R. Alur and P. Madhusudan. Visibly Pushdown Languages. In ACM Symposium on Theory of Computing (STOC), pages 202–211, 2004.
- M. Arenas, L. Bertossi, and J. Chomicki. Consistent Query Answers in Inconsistent Databases. In ACM Symposium on Principles of Database Systems (PODS), pages 68–79, 1999.
- M. Arenas and L. Libkin. XML Data Exchange: Consistency and Query Answering. In ACM Symposium on Principles of Database Systems (PODS), pages 13–24, 2005.
- A Balmin, Y. Papakonstantinou, and V. Vianu. Incremental Validation of XML Documents. ACM Transactions on Database Systems (TODS), 29(4):710–751, December 2004.
- P. Bille. Tree Edit Distance, Alignment and Inclusion. Technical Report TR-2003-23, The IT University of Copenhagen, 2003.
- M. Bojanczyk, C. David, A. Muscholl, T. Schwentick, and L. Segoufin. Twovariable logic on data trees and xml reasoning. In *PODS*, pages 10–19, 2006.
- H. Comon, M. Dauchet, R. Gilleron, C. Löding, F. Jacquemard, D. Lugiez, S. Tison, and M. Tommasi. Tree Automata Techniques and Applications. Available on: http://www.grappa.univ-lille3.fr/tata, 1997. release 2007.
- M. R. Fellows. Parameterized Complexity: The Main Ideas and Connections to Practical Computing. *Electronic Notes in Theoretical Computer Science (ENTCS)*, 61, 2002.
- S. Flesca, F. Furfaro, S. Greco, and E. Zumpano. Repairs and Consistent Answers for XML Data with Functional Dependencies. In *International XML Database* Symposium (Xsym), volume 2824 of Lecture Notes in Computer Science, pages 238–253. Springer, 2003.
- S. Flesca, F Furfaro, S. Greco, and E. Zumpano. Querying and Repairing Inconsistent XML Data. In Web Information Systems Engineering (WISE), pages 175–188, 2005.
- M. Frick, M. Grohe, and C. Koch. Query Evaluation on Compressed Trees. In Logic in Computer Science (LICS), pages 188–200, 2003.
- O. Gauwin, A.-C. Caron, J. Niehren, and S. Tison. Complexity of Earliest Query Answering with Streaming Tree Automata. In *Programming Language Technologies* for XML (PLAN-X), 2008.
- V. Kumar, P. Madhusudan, and M. Viswanathan. Visibly Pushdown Automata for Streaming XML. In International Conference on World Wide Web (WWW), pages 1053–1062, 2007.
- 15. L. Libkin. Data Exchange and Incomplete Information. In ACM Symposium on Principles of Database Systems (PODS), pages 60–69, 2006.
- M. Lohrey and S. Maneth. The Complexity of Tree Automata and XPath on Grammar-compressed Trees. *Theoretical Computer Science (TCS)*, 363(2):196– 210, 2006.
- 17. S. Staworko. Declarative Inconsistencies Handling in Relational and Semistructured Databases. PhD thesis, State University of New York at Buffalo, 2007.
- S. Staworko and J. Chomicki. Validity-Sensitive Querying of XML Databases. In EDBT Workshops (dataX), pages 164–177. Springer, 2006.

- S. Staworko, E. Filiot, and J. Chomicki. Querying Regular Sets of XML Documents. In International Workshop on Logic in Databases (LiD), 2008. Full version available at http://www.grappa.univ-lille3.fr/~staworko/papers/lid08.pdf.
- J. W. Thatcher and Wrightm J. B. Generalized Finite Automata with an Application to a Decision Problem of Second-order Logic. *Mathematical System Theory*, 2:57–82, 1968.
- A. Thomo, S. Venkatesh, and Y. Y. Ye. Visibly Pushdown Transducers for Approximate Validation of Streaming XML. In *Foundations of Information and Knowledge* Systems (FOIKS), pages 219–238, 2008.
- M. Y. Vardi. The Complexity of Relational Query Languages. In ACM Symposium on Theory of Computing (STOC), pages 137–146, 1982.