

# Query Evaluation on Compressed Trees\*

Markus Frick<sup>†</sup>

Martin Grohe<sup>‡</sup>

Christoph Koch<sup>§</sup>

## Abstract

*This article studies the problem of evaluating unary (or node-selecting) queries on unranked trees compressed in a natural structure-preserving way, by the sharing of common subtrees. The motivation to study unary queries on unranked trees comes from the database field, where querying XML documents, which can be considered as unranked labelled trees, is an important task.*

*We give complexity results for the evaluation of XPath and monadic datalog queries. Furthermore, we propose a new automata-theoretic formalism for querying trees and give algorithms for evaluating queries defined by such automata.*

## 1. Introduction

Semi-structured data, best-known in the syntax of XML, have caused a significant paradigm shift in the field of database systems, and have also been one of the central research topics in database theory over the last five years (see [1] and [24] for surveys). While classical relational databases can be described as relational structures, XML-documents are best modelled by unranked trees. This article studies the problem of evaluating unary (or node-selecting) queries on unranked trees compressed in a natural structure-preserving way, by the *sharing of common subtrees*. Node-selecting queries are not only of interest as basic queries in their own right, but are also an important building block for more complex queries. In particular, the node-selecting path query language *XPath* is at the core of several major XML-related technologies, such as XML Query, XML Schema, and XSLT, the principal query language, schema definition formalism, and stylesheet language for XML, respectively. Thus, the efficient processing of XPath queries (and the study of node-selecting XML queries in general) is

paramount to the overall success of all these technologies.

The compression of XML-trees into directed acyclic graphs by the sharing of subtrees has recently been proposed by Buneman, Grohe and Koch [6]. It can be seen as a direct generalisation of the compression of Boolean functions into OBDDs (cf. [5]) used so successfully in symbolic model checking [7, 8]. The approach bears the promise of advancing the state of the art in XML query processing at two fronts. First, compression allows to keep larger document trees in main memory (where they can be efficiently evaluated) and permits a substantial speedup by often avoiding the need to use slow secondary storage when trees are large and otherwise cannot be kept in main memory as a whole. Second, evaluating queries on compressed trees in practice saves time by avoiding redundant computations. In [6], an evaluator was implemented for a large fragment of XPath (*Core XPath*, which was introduced in [17]) and extensively benchmarked on practical XML documents several hundreds of Megabytes large and consisting of trees comprising tens of millions of nodes. The compression ratios obtained were very promising, and the actual efficiency of query processing obtained was astonishing.

The main objective of this article is to create a solid theoretical foundation for the approach. The problem of evaluating Core XPath queries on compressed instances has been shown to be fixed-parameter tractable in [6]. More specifically, the problem has been shown to be solvable in time  $O(k \cdot 2^k \cdot n)$ , where  $k$  denotes the size of the query and  $n$  the size of the compressed instance. Complementing this result, here we prove that the problem is PSPACE-complete. Furthermore, we show that the problem of evaluating queries of the positive Core XPATH fragment (i.e., without negation) is NP-complete. Let us remark that the problem of evaluating Core XPath queries on uncompressed trees is known to be in polynomial time (actually, PTIME-complete [17, 18]).

Even though undoubtedly very important in practice, from a theoretical perspective XPath seems to be a very ad-hoc language that leaves a lot to be desired. *Monadic second-order logic (MSO)* on trees, on the other hand, is well-known to have beautiful theoretical properties. In particular, it has well-balanced expressive power in that it is expressive enough for most purposes, but on the other hand still has good algorithmic properties due to its connection

\*This article is a long version of the extended abstract [14] which appeared in Proc. LICS 2003, Ottawa, Canada.

<sup>†</sup>Laboratory for Foundations of Computer Science, University of Edinburgh, Edinburgh EH9 3JZ, Scotland, UK. Email: mfrick@inf.ed.ac.uk

<sup>‡</sup>Lehrstuhl Logik in der Informatik, Institut für Informatik, Humboldt-Universität zu Berlin, Johann von Neumann-Haus, Haus 4 / 4. Etage, Rudower Chaussee 25, 12489 Berlin (Adlershof). Email: grohe@informatik.hu-berlin.de

<sup>§</sup>Institut für Informationssysteme, Abteilung für Datenbanken und Artificial Intelligence, Technische Universität Wien, Favoritenstraße 9, A-1040 Wien, Austria. Email: koch@dbai.tuwien.ac.at

with tree automata. Indeed, MSO has been proposed as a “benchmark” for the expressive power of node-selecting XML query languages [25]. Nevertheless, MSO itself is not suitable as a practical query language because it allows to express very complex queries very concisely, which makes the query evaluation problem intractable even on uncompressed trees (cf. [27, 13]). But there are nice languages which have the same expressive power as MSO on trees, but admit much more efficient query evaluation. The modal  $\mu$ -calculus may be seen as an example of such a language (at least on ranked trees). In the context of querying XML, the most promising such language is *monadic datalog*. It has the same expressive power as MSO, but admits query evaluation in time linear in both the size of the datalog program and the size of the tree [16].

We study the evaluation problem for monadic datalog queries on compressed instances. We show that, as for the strictly weaker Core XPath, the problem is PSPACE-complete. Of course the PSPACE-hardness was to be expected, but the containment of the problem in PSPACE may be viewed as mildly surprising.

Next, slightly digressing from the main focus of this article, we discuss the connection between compressed binary and unranked trees. Even though XML-documents are naturally represented as unranked trees, we believe that it may be worthwhile to convert them into binary trees first and then only work with compressed binary trees. We show how such a conversion can be carried out without paying too high a price for it and observe that in certain situations the compressed binary instances may be exponentially smaller than the compressed unranked instances they represent. An additional advantage this may have in practice is that nodes of binary instances can be stored in a fixed size memory segment, whereas nodes of unranked instances may have adjacency lists of unbounded length, which tend to lead to high memory fragmentation. Experimental evidence suggests that the conversion to binary instances is worthwhile in practice and may lead to more memory- and time-efficient XML query engines.

Returning to the query evaluation problem, an alternative approach to querying trees can be based on tree-automata.

Right from the beginning, automata-theoretic ideas have played a central role in XML-related research. Automata theory has been used for evaluating path and pattern queries [4, 25, 26, 16], as a basis for XML schema languages [21], for defining XML transducers [22], and for XML data stream processing [19]; See [24] for a survey of automata-theoretic work related to XML.

Most important in our context are *query automata*, proposed by Neven and Schwentick [26] to describe node-selecting queries on unranked trees. We suggest a similar automata model that we call *selecting tree automaton* (STA). STAs have the same querying power and similar al-

gorithmic properties as query automata, but we feel they are much cleaner and simpler. Even though our main interest here is in querying compressed instances, STAs are relevant in the uncompressed setting as well. We show that STA-queries on compressed instances can be evaluated in time  $2^{O(s)} \cdot n$ , where  $s$  is the size of the automaton and  $n$  the size of the compressed instance. Unfortunately, this seems to be too inefficient for practical purposes, because we usually cannot expect our automata to be so small that a factor in the running time that is exponential in  $s$  is acceptable. For example, translating a monadic datalog program into an equivalent STA causes an exponential blow-up in size. Therefore, we also consider a restricted model of *weak selecting tree automata*, whose definition captures the intuitions of monotonic inference in monadic datalog. This is witnessed by the fact that the above transformation from monadic datalog (which comes with the mentioned exponential blow-up) naturally yields weak STAs, and the time to evaluate such automata on compressed instances is linear.

The structure of this article basically follows the order of contributions given above.

## 2. Compressed Trees

In this section, we review the framework for querying trees compressed by subtree sharing that has been introduced in [6]. We emphasise the central role of the familiar notion of bisimilarity in this framework.

**2.1. Instances and Bisimilarity.** A *schema* is a finite set of unary relation names. Let  $\sigma = \{S_1, \dots, S_n\}$  be a schema. An *instance of schema*  $\sigma$ , or  $\sigma$ -*instance*, is a tuple

$$\mathbf{I} = (V^{\mathbf{I}}, \gamma^{\mathbf{I}}, \text{root}^{\mathbf{I}}, S_1^{\mathbf{I}}, \dots, S_n^{\mathbf{I}}),$$

where  $V^{\mathbf{I}}$  is the (finite) set of vertices,  $\gamma^{\mathbf{I}} : V^{\mathbf{I}} \rightarrow (V^{\mathbf{I}})^*$  is a function whose graph is acyclic and has a unique vertex of in-degree 0 from which all other vertices are reachable,  $\text{root}^{\mathbf{I}}$  is this vertex of in-degree 0, and  $S_1^{\mathbf{I}}, \dots, S_n^{\mathbf{I}}$  are subsets of  $V^{\mathbf{I}}$ .

Here the *graph of*  $\gamma$  is the directed graph with vertex set  $V^{\mathbf{I}}$  and an edge from  $v$  to  $w$  if  $w$  occurs in  $\gamma(v)$ . We call this graph the *DAG of*  $\mathbf{I}$ . Occasionally we denote its edge relation by  $E^{\mathbf{I}}$ . For vertices  $u, v \in V$  we write  $u \xrightarrow{i} v$  if there is an  $n \geq i$  and vertices  $v_1, \dots, v_n$  such that  $\gamma(u) = v_1 \dots v_n$  and  $v = v_i$ . Intuitively, if  $u \xrightarrow{i} v$  then  $v$  is the *i-th child* of  $u$  and  $u$  is the *parent* of  $v$ . An instance is *k-ary*, for some  $k \geq 1$ , if every vertex has at most  $k$  children. If the instance  $\mathbf{I}$  is clear from the context, we often omit the superscript  $\mathbf{I}$ .

A *tree instance* is an instance whose DAG is a tree. Tree-instances are our model of XML-documents. The unary relations represented by the schema are used to encode the

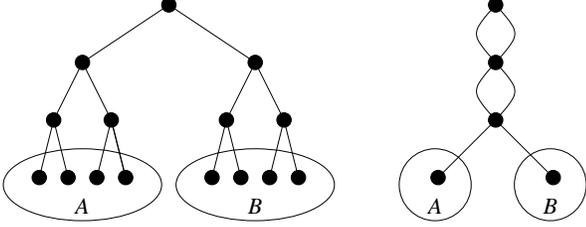


Figure 1. Two bisimilar instances over schema  $\sigma = \{A, B\}$ .

relevant information carried by the vertices of an XML-tree. This may be the XML-tag of a vertex, but also information encoded in the alphanumeric data at the vertex. Indeed, the latter is of crucial importance because node-selecting queries will usually access the alphanumeric data. The implementation of Core XPath in [6] admits querying alphanumeric data through so-called *string constraints*.

A *bisimilarity relation* between two  $\sigma$ -instances  $\mathbf{I}$  and  $\mathbf{J}$  is a binary relation  $\sim \subseteq V^{\mathbf{I}} \times V^{\mathbf{J}}$  such that for all  $v \in V^{\mathbf{I}}, w \in V^{\mathbf{J}}$  with  $v \sim w$  we have

- for all  $i$  and  $v' \in V^{\mathbf{I}}$ , if  $v \xrightarrow{i} v'$  then there exists  $w' \in V^{\mathbf{J}}$  such that  $w \xrightarrow{i} w'$  and  $v' \sim w'$ ,
- for all  $i$  and  $w' \in V^{\mathbf{J}}$ , if  $w \xrightarrow{i} w'$  then there exists  $v' \in V^{\mathbf{I}}$  such that  $v \xrightarrow{i} v'$  and  $v' \sim w'$ , and
- for all  $S \in \sigma$ : ( $v \in S^{\mathbf{I}} \iff w \in S^{\mathbf{J}}$ ).

Thus, our notion of bisimilarity relation takes both node and edge labels that are present in the tree into account.

If there is some bisimilarity relation  $\sim$  between  $\mathbf{I}$  and  $\mathbf{J}$  such that  $v \sim w$ , we call the vertices  $v$  and  $w$  *bisimilar* (we write  $v \approx w$  or  $(\mathbf{I}, v) \approx (\mathbf{J}, w)$ ). The instances  $\mathbf{I}$  and  $\mathbf{J}$  are *bisimilar* (we write  $\mathbf{I} \approx \mathbf{J}$ ) if  $root^{\mathbf{I}} \approx root^{\mathbf{J}}$ . An example of two bisimilar instances is given in Figure 1.

If  $\mathbf{I}$  is an instance and  $\sim$  a bisimilarity relation on  $\mathbf{I}$  (that is, between  $\mathbf{I}$  and  $\mathbf{I}$ ), then  $\mathbf{I}/\sim$  is the instance obtained from  $\mathbf{I}$  by identifying all vertices  $v, w$  with  $v \sim w$ . Pairs of edges  $v \xrightarrow{i} w, v' \xrightarrow{j} w'$  (with  $v \sim v', w \sim w'$ ) may only be identified in this process if  $i = j$ . We define a partial order  $\preceq$  on the class of all  $\sigma$ -instances by letting  $\mathbf{I} \preceq \mathbf{J}$  if there is a bisimilarity relation  $\sim$  on  $\mathbf{J}$  such that  $\mathbf{I}$  is isomorphic to  $\mathbf{J}/\sim$ . More precisely,  $\preceq$  is a partial order on the set of all isomorphism classes of instances. But no harm is done by blurring this distinction here and in the following.

**Lemma 1 ([6]).** *Let  $\mathbf{I}$  be a  $\sigma$ -instance and let  $\mathcal{L}(\mathbf{I})$  be the class of all instances bisimilar to  $\mathbf{I}$ . Then  $(\mathcal{L}(\mathbf{I}), \preceq)$  is a lattice. Its maximal element  $\mathbf{T}(\mathbf{I})$  is the only tree instance in  $\mathcal{L}(\mathbf{I})$ . The minimal element  $\mathbf{M}(\mathbf{I})$  is also characterised by the fact that it contains the least number of vertices of all instances in  $\mathcal{L}(\mathbf{I})$ .*

It will be necessary later to have a canonical definition of  $\mathbf{T}(\mathbf{I})$ , and not just a characterisation up to isomorphism. If  $v, v'$  are vertices in an instance  $\mathbf{I}$ , and there are intermediate vertices  $v_1, \dots, v_{n-1}$  such that  $v \xrightarrow{i_1} v_1 \dots v_{n-1} \xrightarrow{i_n} v'$ , we say that the integer sequence  $i_1 \dots i_n$  is an *edge-path* between  $v$  and  $v'$ . For each vertex  $v \in V$  we define

$$\Pi(v) = \{P \mid P \text{ is an edge-path from } root \text{ to } v\},$$

and for a set  $S \subseteq V$  we let  $\Pi(S) = \bigcup_{v \in S} \Pi(v)$ . Note that the vertices of  $\mathbf{T}(\mathbf{I})$  are in one-to-one correspondence with the elements of  $\Pi(V)$ . Thus we can define our canonical representation of  $\mathbf{T}(\mathbf{I})$  to have vertex set  $\Pi(V)$  (and all relations defined in the obvious way). Also note that for bisimilar instances  $\mathbf{I} \approx \mathbf{J}$  and vertices  $v \in V^{\mathbf{I}}, w \in V^{\mathbf{J}}$  we have  $v \approx w$  if and only if  $\Pi(v) = \Pi(w)$ .

**2.2. Representations and Size.** Our machine model is a standard RAM model with addition and subtraction as arithmetic operations. We use a uniform cost measure. While for some of the theoretical considerations of this article a logarithmic cost measure would be nicer (cf. Remark 12), we think that for the practical analysis of our algorithms a uniform measure is most appropriate. After all, a main motivation for this research is to give main memory algorithms for querying XML-documents, and this basically means that in our algorithms we never have to handle numbers that do not fit into a single memory word.

We represent instances in a straightforward way based on an adjacency list representation of the underlying DAG. There is one important twist: Instances may contain multiple edges, and instead of storing them all separately, we just use one adjacency list entry which contains an integer representing the multiplicity to represent consecutive edges from a vertex to a child. Since the order of the children of a vertex is important, we can only do this for consecutive edges to the same child (see Figure 2). In practice, this concise representation of multiple edges is extremely important, because XML-trees tend to be very broad and shallow, and therefore their compressed versions tend to have many multiple edges. We denote the size of the representation of an instance  $\mathbf{I}$  by  $\|\mathbf{I}\|$ . Note that we have  $|V^{\mathbf{I}}| \leq O(\|\mathbf{I}\|)$  and  $\|\mathbf{I}\| \leq O(|V^{\mathbf{I}}|^2)$ . Moreover, we have  $\|\mathbf{I}\| \leq O(|E^{\mathbf{I}}|)$ , but our concise representation of multiple edges and the uniform cost model imply that we cannot bound  $|E^{\mathbf{I}}|$  in terms of  $\|\mathbf{I}\|$ . If we let  $\text{mult}(\mathbf{I})$  be the maximum number of consecutive multiple edges, that is, the maximum number of edges represented by a single adjacency list entry, we have  $|E^{\mathbf{I}}| \leq O(\|\mathbf{I}\| \cdot \text{mult}(\mathbf{I}))$ .

The following earlier result refers to the computation of compressed instances even *with edge multiplicities*:

**Theorem 2 ([6]).** *There is an algorithm that, given an instance  $\mathbf{I}$ , computes the minimal bisimilar instance  $\mathbf{M}(\mathbf{I})$  in time  $O(\|\mathbf{I}\|)$ .*



Figure 2.

### 3. Queries and Query Languages

**3.1. Queries.** Since we think of an instance  $I$  as being a compressed representation of the tree instance  $T(I)$ , we define the semantics of queries with respect to the tree instances. Our notion of query is based on the use of this term in finite model theory: A (unary) query  $Q$  of schema  $\sigma$  associates with every tree instance  $T$  of schema  $\sigma$  a subset  $Q(T) \subseteq V^T$  in such a way that for every isomorphism  $\pi$  from a tree-instance  $T$  to a tree instance  $T'$  we have  $\pi(Q(T)) = Q(T')$ . Note that the term “query” usually has a different meaning in the theory of semi-structured data; there a query is a mapping from tree instances to tree instances (see [1]). What we call unary query here is usually called *pattern* in this framework.

We want to evaluate queries on compressed instances, preferably without fully decompressing them. The problem is that we cannot always represent the result of a query in the instance we are given: While every subset  $X \subseteq V^I$  canonically corresponds to the subset  $\Pi(X) \subseteq V^{T(I)} = \Pi(V^I)$ , unless  $I = T(I)$  it is not the case that for every set  $Y \subseteq V^{T(I)}$  there is a set  $X \subseteq V^I$  such that  $Y = \Pi(X)$ . So to represent the answer of a query we may have to partially decompress the instance. The following definition makes this precise:

**Definition 3.** The *evaluation problem* for a query language  $L$  on a class  $C$  of instances is the following problem:

*Input:* Instance  $I \in C$  and query  $Q \in L$ .  
*Problem:* Return an instance  $J$  and a subset  $\tilde{Q} \subseteq V^J$  such that  $I$  and  $J$  are bisimilar and  $Q(T(I)) = \Pi(\tilde{Q})$ .

If  $C$  is not explicitly mentioned, it is understood to be the class of all instances.

Our complexity-theoretic results only refer to the *decision version* of the evaluation problem:

*Input:* Instance  $I \in C$ , vertex  $v \in V^{T(I)}$ , and query  $Q \in L$ .  
*Problem:* Decide if  $v \in Q(T(I))$ .

In this definition, we supply a node of the *uncompressed* tree-version of the instance  $I$  with the input, since nodes

of  $I$  do not necessarily exist in  $Q(I)$ . Clearly, a good way to formulate global properties of instances as decision problems is to check queries on the root node, since  $\Pi(\text{root}^I) = \{\text{root}^{T(I)}\}$  for all instances.

**3.2. Complexity.** We assume that the reader is familiar with the standard complexity classes such as PTIME, NP, and PSPACE. It is convenient to phrase some of our results in the terminology of fixed-parameter tractability (see [10], or [20] for a short introduction into the notions most relevant here). Actually, we only need one definition: The evaluation problem for  $L$  on  $C$  is *fixed parameter tractable* if there is a computable function  $f$ , a constant  $c$ , and an algorithm solving the problem in time  $f(k) \cdot n^c$ , where  $n$  is the size of the input instance and  $k$  the size of the input query.

**3.3. Logic and relational structures.** We assume that the reader is familiar with relational structures, first-order logic FO, and monadic second-order logic MSO (see, for example, [11]). In the logical context, we describe  $\sigma$ -tree instances as relational structures whose vocabulary consists of all unary relation symbols in the schema  $\sigma$  and, in addition, the unary relation symbols *Root*, *Leaf*, *Last-Sibling*, and the binary relation symbols *First-Child* and *Next-Sibling*, all with the natural meanings (cf. [16]). We occasionally call *Root*, *Leaf*, *Last-Sibling*, *First-Child*, and *Next-Sibling* the *built-in predicates*. We use  $T$  to denote both the tree instance and the relational structure representing it.

If  $T$  is a tree instance and  $\varphi(x)$  an MSO-formula with one free variable, then we let  $\varphi(T)$  be the set of all vertices  $v \in V^T$  such that  $T$  satisfies  $\varphi(x)$  if  $x$  is interpreted by  $v$ . We call  $T \mapsto \varphi(T)$  the query defined by  $\varphi$ .

**3.4. Monadic datalog.** We assume that the reader is familiar with datalog, which may be viewed as logic programming without function symbols (cf. [2]). A datalog program is *monadic* if all its IDB predicates (that is, intensional predicates that appear in rule heads somewhere in the program) are unary. We interpret monadic datalog programs over tree instances. A monadic datalog program of schema  $\sigma$  may use as EDB predicates (that is, extensional predicates which are determined by the structure the program is interpreted over) the built-in predicates *Root*, *Leaf*, *Last-Sibling*, the binary relation symbols *First-Child* and *Next-Sibling*, the predicates in  $\sigma$ , and a predicate  $\bar{S}$  for every  $S \in \sigma$  which is interpreted as the complement of  $S$ . Each program  $\mathcal{P}$  has a distinguished *goal (IDB) predicate*. The query defined by  $\mathcal{P}$  maps a tree instance  $T$  to the set of all vertices  $v$  such that  $\mathcal{P}$  derives over  $T$  that  $v$  is in the goal predicate. Two programs are *equivalent* if they define the same query.

A *valuation* on a datalog rule  $r$  is a mapping  $\varphi$  which maps each variable of  $r$  to a domain element. Given an

atom  $p(x_1, \dots, x_m)$ , let

$$\varphi(p(x_1, \dots, x_m)) := p(\varphi(x_1), \dots, \varphi(x_m)).$$

The popular fixpoint semantics of (monadic) datalog can be defined by means of a small-step<sup>1</sup> monotonic *immediate consequence operator*  $\mathcal{T}_{\mathcal{P}}$ . Given a set of ground (i.e., variable-free) atoms  $X$ ,  $\mathcal{T}_{\mathcal{P}}$  chooses a rule of  $\mathcal{P}$  (with head predicate  $P$ ) and a node  $v$  such that  $X$  satisfies the rule body and a new ground atom  $P(v)$  can be inferred (then,  $\mathcal{T}_{\mathcal{P}}(X) = X \cup \{P(v)\}$ ). We write the  $k$ -times iterated application of  $\mathcal{T}_{\mathcal{P}}$  as  $\mathcal{T}_{\mathcal{P}}^k$  and the fixpoint  $\mathcal{T}_{\mathcal{P}}^m = \mathcal{T}_{\mathcal{P}}^{m+1}$  as  $\mathcal{T}_{\mathcal{P}}^\omega$ . In this context, we may consider an instance  $\mathbf{T}$  as a set of unary and binary ground atoms and evaluate  $\mathcal{P}$  on  $\mathbf{T}$  as  $\mathcal{T}_{\mathcal{P}}^\omega(\mathbf{T})$ .

We only use monadic datalog programs with a restricted syntax described next. In a *TMNF program* (“Tree-marking Normal Form”), each rule is an instance of one of the four rule templates (with “types” 1 to 4)

$$\begin{aligned} P(x) &\leftarrow U(x). & (1) \\ P(x) &\leftarrow P_0(x_0) \wedge B(x_0, x). & (2) \\ P(x_0) &\leftarrow P_0(x) \wedge B(x_0, x). & (3) \\ P(x) &\leftarrow P_1(x) \wedge P_2(x). & (4) \end{aligned}$$

where  $P, P_0, P_1, P_2$  are IDB predicates and  $U, B$  are EDB predicates.

**Proposition 4 ([16]).** *Every monadic datalog program (over trees) can be translated into an equivalent TMNF program in linear time.*

Thus, in the following, we only deal with programs in TMNF. All worst-case bounds for TMNF query evaluation translate immediately to the evaluation of monadic datalog. Note that monadic datalog captures MSO over trees [16] and a program  $\mathcal{P}$  can be evaluated in time  $O(|\mathbf{T}| * |\mathcal{P}|)$  when  $\mathbf{T}$  is a tree-instance [16].

**3.5. Core XPath.** XPath uses thirteen binary relations – called axes – for navigating in trees [30]. We only need to introduce five of them in this article, *Self* (the identity relation on  $V$ ), *Child* (the intuitive child relation;  $Child(v, w)$  iff  $w$  is a child of  $v$ ), *Parent* (its inverse), *Descendant* (the transitive closure of *Child*), and *Ancestor* (its inverse). Given a binary axis relation  $\chi$ , an axis is its *inverse*, denoted  $\chi^{-1}$ , if  $\langle x, y \rangle \in \chi \Leftrightarrow \langle y, x \rangle \in \chi^{-1}$ . For each XPath axis, there is also an XPath axis that is its inverse. For example,  $Self^{-1} = Self$ ,  $Child^{-1} = Parent$ , and  $Descendant^{-1} = Ancestor$ . In the following definition of a

<sup>1</sup>Note that this nondeterministic definition is slightly nonstandard but makes certain induction proofs more straightforward. Usually, the immediate consequence operator adds to  $X$  in one step *all* the atoms that can be inferred as described here from the atoms in  $X$ . Of course, the fixpoint of the both versions of the operator is the same on all instances and programs.

fragment of XPath, we assume all 13 axes to be supported, even if only a few have been introduced here (for a complete formal definition of Core XPath see [17]).

**Definition 5.** Let  $\mathbf{T}$  be a tree-instance. We define the syntax of *Core XPath* by the EBNF

$$\begin{aligned} \text{corexpath:} & \quad \text{locationpath} \mid \text{'/' locationpath} \\ \text{locationpath:} & \quad \text{locationstep ('/' locationstep)*} \\ \text{locationstep:} & \quad \chi \text{'::'} P \mid \chi \text{'::'} P \text{'[' pred ']} \\ \text{pred:} & \quad \text{pred 'and' pred} \mid \text{pred 'or' pred} \\ & \quad \mid \text{'not(' pred ')'} \mid \text{corexpath} \mid \text{'(' pred ')'} \end{aligned}$$

“corexpath” is the start production,  $\chi$  stands for an axis, and  $P$  for a “node test”, that is, a unary relation from  $\sigma$  or “\*”, meaning “any node”  $V^{\mathbf{T}}$ .

The semantics of Core XPath queries on tree-instances  $\mathbf{T}$  is defined by two functions  $\mathcal{S}$  and  $\mathcal{E}$  (for Core XPath expressions and condition predicates, respectively):

$$\begin{aligned} \mathcal{S} &: \mathcal{L}(\text{corexpath}) \rightarrow 2^{V^{\mathbf{T}} \times V^{\mathbf{T}}} \\ \mathcal{S}[\chi::P[e]] &:= \{\langle x, y \rangle \mid \chi^{\mathbf{T}}(x, y) \wedge y \in (P \cap \mathcal{E}[e])\} \\ \mathcal{S}[/\pi] &:= V^{\mathbf{T}} \times \{x \mid \langle \text{root}^{\mathbf{T}}, x \rangle \in \mathcal{S}[\pi]\} \\ \mathcal{S}[\pi_1/\pi_2] &:= \{\langle x, z \rangle \mid \exists y : \langle x, y \rangle \in \mathcal{S}[\pi_1] \wedge \langle y, z \rangle \in \mathcal{S}[\pi_2]\} \end{aligned}$$

$$\begin{aligned} \mathcal{E} &: \mathcal{L}(\text{pred}) \rightarrow 2^{V^{\mathbf{T}}} \\ \mathcal{E}[e_1 \text{ and } e_2] &:= \mathcal{E}[e_1] \cap \mathcal{E}[e_2] \\ \mathcal{E}[e_1 \text{ or } e_2] &:= \mathcal{E}[e_1] \cup \mathcal{E}[e_2] \\ \mathcal{E}[\text{not}(e)] &:= V^{\mathbf{T}} - \mathcal{E}[e] \\ \mathcal{E}[\pi] &:= \{x_0 \mid \exists x : \langle x_0, x \rangle \in \mathcal{S}[\pi]\} \end{aligned}$$

Here,  $\pi, \pi_1$  and  $\pi_2$  are location paths. Query  $Q$  results in the set  $\{y \mid \exists x : \langle x, y \rangle \in \mathcal{S}[Q]\}$ .

An example of a Core XPath query is

$$\text{/descendant::*[child::A and child::B]/child::*},$$

which selects all children of descendants of the root node that (i.e., the descendants) have a child node labelled A and a child node labeled B.

Core XPath is a strict fragment of XPath [30], both syntactically and semantically.

A mapping from Core XPath to monadic datalog with stratified negation was given in [15]; we strengthen this result to (negation-free) TMNF. Note again that we assume that Core XPath supports all XPath axes.

**Proposition 6.** *Every Core XPath query can be translated into an equivalent TMNF program in linear time and logarithmic space.*

$$\begin{aligned}
Child_P(x) &\leftarrow P(x_0) \wedge First-Child(x_0, x). \\
Child_P(x) &\leftarrow Child_P(x_0) \wedge Next-Sibling(x_0, x). \\
\overline{Child_P}(x) &\leftarrow Root(x). \\
\overline{Child_P}(x) &\leftarrow \overline{P}(x_0) \wedge First-Child(x_0, x). \\
\overline{Child_P}(x) &\leftarrow \overline{Child_P}(x_0) \wedge Next-Sibling(x_0, x).
\end{aligned}$$

$$\begin{aligned}
Parent_P(x) &\leftarrow Aux-Parent_P(x_0) \wedge \\
&\quad First-Child(x, x_0). \\
Aux-Parent_P(x) &\leftarrow P(x). \\
Aux-Parent_P(x) &\leftarrow Aux-Parent_P(x_0) \wedge \\
&\quad Next-Sibling(x, x_0). \\
\overline{Parent_P}(x) &\leftarrow Leaf(x). \\
\overline{Parent_P}(x_0) &\leftarrow First-Child(x_0, x) \wedge \\
&\quad Aux-Parent_P(x). \\
Aux-Parent_P(x) &\leftarrow \overline{P}(x) \wedge Last-Sibling(x). \\
Aux-Parent_P(x) &\leftarrow \overline{P}(x) \wedge Next-Sibling(x, y) \wedge \\
&\quad Aux-Parent_P(y).
\end{aligned}$$

Figure 3. Axis encodings (for *Child*, *Parent*, and their complements).

*Proof of Proposition 6:* To keep the encoding of this proof simple, we do not strictly adhere to TMNF syntax. However, all datalog rules can be transformed into TMNF by folding pairs of unary atoms from long rules into separate rules. For instance, a rule

$$P(x) \leftarrow P_1(x) \wedge P_2(x) \wedge P_3(x).$$

rewrites into

$$\begin{aligned}
P(x) &\leftarrow P'(x) \wedge P_3(x). \\
P'(x) &\leftarrow P_1(x) \wedge P_2(x).
\end{aligned}$$

where  $P'$  is a new predicate. This is not an additional transformation that has to be composed with the proof, but can be applied a priori to our encoding.

The encoding presented here is closely based on one presented in [15]; however, there, negation was encoded using datalog with stratified negation, which we do not want to use here.<sup>2</sup> The idea is to push down negation to unary built-in predicates, for which we, by the definitions of Section 3.4, have the complements available. The slightly

<sup>2</sup>Even though extending monadic datalog with stratified negation changes neither the complexity nor the expressiveness of the formalism.

tricky part is the (binary) axis relations, which are also the only part of a Core XPath query for which recursion is needed in our encoding. As in [15], for each axis  $\chi$ , we can set up a basically fixed template program (modulo occurrences of IDB predicate  $P$ ) defining a predicate  $\chi_P$  with the intuitive meaning

$$\chi_P(x) \leftarrow P(x_0) \wedge \chi(x_0, x).$$

As TMNF with the given built-in relations provides all the machinery to check a given predicate  $P$  on *each* of the nodes of a region in the tree relevant to an axis, negated axes  $\overline{\chi_P}$  (denoting the complement of the set of nodes to which  $\chi_P$  evaluates) can be encoded as well. We show implementations of  $Child_P$ ,  $Parent_P$ ,  $\overline{Child_P}$ , and  $\overline{Parent_P}$  in Figure 3 and an implementation of  $Ancestor_P$  in the example of Figure 4.

First we provide the encoding for positive Core XPath, i.e. for Core XPath without negation. This is precisely as in [15]. Given a query  $\pi$ , we obtain a program computing  $\pi$  by starting with the query predicate  $\mathcal{S}_R[\pi]$  and including rules according to the patterns shown below until all IDB predicates have been defined.

$$\begin{aligned}
\mathcal{S}_R[\chi::P[e]](x) &\leftarrow \chi_V(x) \wedge P(x) \wedge \mathcal{E}[e](x). \\
\mathcal{S}_R[/\chi::P[e]](x) &\leftarrow \chi_{Root}(x) \wedge P(x) \wedge \mathcal{E}[e](x). \\
\mathcal{S}_R[\pi/\chi::P[e]](x) &\leftarrow \chi_{\mathcal{S}_R[\pi]}(x) \wedge P(x) \wedge \mathcal{E}[e](x). \\
Aux-\mathcal{S}_L[P[e]](x) &\leftarrow P(x) \wedge \mathcal{E}[e](x). \\
\mathcal{S}_L[\chi::P[e]](x) &\leftarrow \chi_{Aux-\mathcal{S}_L[P[e]]}^{-1}(x). \\
Aux-\mathcal{S}_L[P[e]/\pi](x) &\leftarrow \mathcal{S}_L[\pi](x) \wedge P(x) \wedge \mathcal{E}[e](x). \\
\mathcal{S}_L[\chi::P[e]/\pi](x) &\leftarrow \chi_{Aux-\mathcal{S}_L[P[e]/\pi]}^{-1}(x). \\
\mathcal{E}[\pi](x) &\leftarrow \mathcal{S}_L[\pi](x). \\
\mathcal{E}[e_1 \text{ and } e_2](x) &\leftarrow \mathcal{E}[e_1](x) \wedge \mathcal{E}[e_2](x).
\end{aligned}$$

$$\mathcal{E}[e_1 \text{ or } e_2](x) \leftarrow \mathcal{E}[e_1](x). \quad \mathcal{E}[e_1 \text{ or } e_2](x) \leftarrow \mathcal{E}[e_2](x).$$

Note that as in the encoding of [15], query subexpressions within brackets (i.e., conditions) are intuitively “reversed” (using the  $\mathcal{S}_L$  predicates and the axis inverses  $\chi^{-1}$ ) to direct all computation in the query tree towards the “hot point” that selects nodes.

Now for negation, inside boolean condition expressions (with “and”, “or” and “not” as operations and path expressions assumed atomic), negations are pushed down as far as possible using De Morgan’s laws (which of course do not lead to an increase in formula size). Then, we proceed as above with the additional rule templates

$$\begin{aligned}
\mathcal{E}[\text{not}(\pi)](x) &\leftarrow \overline{\mathcal{S}_L[\pi]}(x). \\
\overline{\mathcal{S}_L[\chi::P[e]]}(x) &\leftarrow \chi_{Aux-\mathcal{S}_L[P[e]]}^{-1}(x). \\
\overline{\mathcal{S}_L[\chi::P[e]/\pi]}(x) &\leftarrow \chi_{Aux-\mathcal{S}_L[P[e]/\pi]}^{-1}(x).
\end{aligned}$$

Rules defining predicates such as  $\overline{\text{Aux-}\mathcal{S}_L[P[e]/\pi]}$  can be easily obtained by again applying De Morgan's law, by which the negation of a rule

$$P(x) \leftarrow P_1(x) \wedge \dots \wedge P_n(x).$$

becomes

$$\overline{P}(x) \leftarrow \overline{P_1}(x). \quad \dots \quad \overline{P}(x) \leftarrow \overline{P_n}(x).$$

Even if not obvious from this rewriting-based discussion, this translation can be carried out in linear time and logarithmic space. Figure 4 gives an example.  $\square$

A natural restriction on Core XPath is to exclude negation. We call the language obtained *positive Core XPath*. This fragment is also interesting as while both monadic datalog and Core XPath are P-complete w.r.t. combined complexity on tree-instances, positive Core XPath is LOGCFL-complete and thus effectively parallelizable [18].

#### 4. Complexity and Evaluation of Monadic Datalog

In this section, we study the complexity of query evaluation for TMNF and Core XPath on compressed instances. Edge multiplicities, as introduced before, may lead to some difficulty; we assume such multiplicities not to be present in the instances of this section. Since they can have practical impact on the compression rate (cf. [6]), we refer to Section 5 where a mapping of unranked trees with edge multiplicities to binary trees is described.

Before we arrive at our main complexity results, we provide a polynomial-space algorithm that computes the predicates derivable for a given node, compressed instance, and TMNF program.

Let  $\mathbf{T}$  be a tree-instance and  $v \in V^{\mathbf{T}}$  a node. Moreover, let  $\mathbf{T}_v$  denote the subtree of  $\mathbf{T}$  rooted by  $v$ .  $\mathbf{T}_v$  is again an instance. Let  $X$  be a set of ground atoms and  $V$  a set of nodes. Then, by  $X/V$  we denote  $\{P(v) \in X \mid v \in V, P \in \text{IDB}(\mathcal{P})\}$ ,  $X/v$  is short for  $X/\{v\}$ , and we abbreviate  $X/V^{\mathbf{T}}$  as  $X/\mathbf{T}$ .

In analogy to  $\mathbf{T}_v$  (the subtree of  $\mathbf{T}$  rooted by  $v$ ), let  $\overline{\mathbf{T}}_v$  denote the complement (or *envelope*) of  $\mathbf{T}_v$ , the tree obtained by removing all of  $\mathbf{T}_v$  from  $\mathbf{T}$  except for  $v$  (which becomes a leaf). Just like  $\mathbf{T}_v$ ,  $\overline{\mathbf{T}}_v$  is again a tree-instance.

The following central lemma now states a stronger version of the fact that (a) in order to decide whether an atom inside a subtree  $\mathbf{T}_v$  can be derived with program  $\mathcal{P}$  (i.e. whether the atom is in  $\mathcal{T}_{\mathcal{P}}^{\omega}(\mathbf{T})$ ), all we need to know is the structure of  $\mathbf{T}_v$  itself and the atoms that are derivable for node  $v$  on  $\mathbf{T}$  (i.e., the atoms of  $\mathcal{T}_{\mathcal{P}}^{\omega}(\mathbf{T})/v$ ), and (b) the converse fact with  $\mathbf{T}_v$  and  $\overline{\mathbf{T}}_v$  exchanged.

The  $\mathcal{T}_{\mathcal{P}}$  operator is nondeterministic, but we may compose a particular run  $\mathcal{T}_{\mathcal{P}}^k(X)$  into a new deterministic operator  $\mathcal{O}_{\mathcal{P}}$ . We denote this by  $\mathcal{O}_{\mathcal{P}} := (\mathcal{T}_{\mathcal{P}}^k)[X]$ . We assume

```

function fixp( $v_0$ : node,
               $X_0$ : set of ground IDB atoms over  $v_0$ ,
               $k$ : integer  $\geq 0$ )
  /* let  $v_0$  have  $n \geq k$  children  $v_1 \dots v_n$  */
  returns set of ground IDB atoms over  $v_k$ 
begin
   $X := X_0 \cup \mathbf{T}/\{v_0, v_1, \dots, v_n\}$ ;
  while no fixpoint reached do
     $X := \mathcal{T}_{\mathcal{P}}(X) \cup$ 
       $\text{fixp}(v_1, X/v_1, 0) \cup \dots \cup \text{fixp}(v_n, X/v_n, 0)$ ;
  return  $X/v_k$ ;
end.

```

Figure 5. The function fixp.

$\mathcal{O}_{\mathcal{P}}$  to choose the same rules and nodes as the particular run of  $\mathcal{T}_{\mathcal{P}}^k$  in the same order. We extend the applicability of  $\mathcal{O}_{\mathcal{P}}$  to sets of ground atoms  $Y \subset X$  by stepping over rules that cannot fire because their bodies are not satisfied in  $Y$  without doing anything.

**Lemma 7.** *Let  $\mathcal{P}$  be a TMNF program,  $\mathbf{T}$  a tree-instance, and  $\mathcal{O}_{\mathcal{P}} := (\mathcal{T}_{\mathcal{P}}^k)[\mathbf{T}]$ . Then,*

$$\mathcal{O}_{\mathcal{P}}(\mathbf{T})/\mathbf{T}_v = \mathcal{O}_{\mathcal{P}}(\mathbf{T}_v \cup \mathcal{O}_{\mathcal{P}}(\mathbf{T})/v)$$

and

$$\mathcal{O}_{\mathcal{P}}(\mathbf{T})/\overline{\mathbf{T}}_v = \mathcal{O}_{\mathcal{P}}(\overline{\mathbf{T}}_v \cup \mathcal{O}_{\mathcal{P}}(\mathbf{T})/v).$$

*Proof of Lemma 7:* By induction on the  $k$  steps of the operator. There are two cases of rules. In the first (which corresponds to TMNF rule templates (1) and (4)), the body contains only a single variable. Thus, the rule application has only strictly local impact, on a single node. In the second case (corresponding to TMNF rule templates (2) and (3)), a rule has at most one single binary atom and a single further IDB atom in the body (which may depend on previously derived facts). Here, a new fact  $P(v)$  may enable the derivation of facts on *different* nodes, but as an immediate consequence, only facts on nodes adjacent in  $\mathbf{T}$  w.r.t. the *First-Child* or *Next-Sibling* relations can be inferred. An atom  $P(v)$  can only contribute to the derivation of another atom  $Q(w)$ , where  $v$  and  $w$  are not adjacent in  $\mathbf{T}$ , if first a trace of new atoms is computed for each of the nodes on the undirected path from  $v$  to  $w$  in  $\mathbf{T}$ . Let  $v'$  be a node on this undirected path between  $v$  and  $w$ , and  $P'(v')$  this witness atom following from  $P(v)$ . Then, subsequently, it suffices to know  $P'(v')$  to infer  $Q(w)$ ;  $P(v)$  is not needed.  $\square$

**Lemma 8.** *Let  $\mathbf{T}$  be a tree-instance,  $v_0 \in V^{\mathbf{T}}$  a node with  $n \geq 0$  children  $v_1 \dots v_n$ ,  $X_0$  a set of ground IDB atoms over node  $v_0$ , and  $0 \leq k \leq n$  an integer. Then, the pseudocode of Figure 5 defines the function*

$$\text{fixp}(v_0, X_0, k) = \mathcal{T}_{\mathcal{P}}^{\omega}(\mathbf{T}_{v_0} \cup X_0)/v_k$$

$$\begin{aligned}
Child_{Root}(x) &\leftarrow Root(x_0) \wedge First-Child(x_0, x). \\
Child_{Root}(x) &\leftarrow Child_{Root}(x_0) \wedge Next-Sibling(x_0, x). \\
\mathcal{S}_R[\text{/child::A[not(descendant::B)]]](x) &\leftarrow Child_{Root}(x) \wedge A(x) \wedge \mathcal{E}[\text{not(descendant::B)}](x). \\
\mathcal{E}[\text{not(descendant::B)}](x) &\leftarrow \overline{\mathcal{S}_L[\text{descendant::B}]}(x). \\
\overline{\mathcal{S}_L[\text{descendant::B}]}(x) &\leftarrow \overline{Ancestor_{Aux-S_L[B]}}(x). \\
\overline{Ancestor_{Aux-S_L[B]}}(x) &\leftarrow Leaf(x). \\
\overline{Ancestor_{Aux-S_L[B]}}(x) &\leftarrow First-Child(x, y) \wedge \overline{Aux-Ancestor_{Aux-S_L[B]}}(y). \\
\overline{Aux-Ancestor_{Aux-S_L[B]}}(x) &\leftarrow \overline{Ancestor_{Aux-S_L[B]}}(x) \wedge \overline{Aux-S_L[B]}(x) \wedge Last-Sibling(x). \\
\overline{Aux-Ancestor_{Aux-S_L[B]}}(x) &\leftarrow \overline{Ancestor_{Aux-S_L[B]}}(x) \wedge \overline{Aux-S_L[B]}(x) \wedge \\
&\quad Next-Sibling(x, y) \wedge \overline{Aux-Ancestor_{Aux-S_L[B]}}(y). \\
\overline{Aux-S_L[B]}(x) &\leftarrow \overline{B}(x).
\end{aligned}$$

Figure 4. Encoding for Core XPath query `/child::A[not(descendant::B)]`.  $A$  and  $B$  (and  $\overline{B}$ ) are predicates in the schema and  $\mathcal{S}_R[\text{/child::A[not(descendant::B)]]]$  is the goal predicate.

which satisfies the two equations

$$\begin{aligned}
\mathcal{T}_{\mathcal{P}}^{\omega}(\mathbf{T})/_{root} &= fi\,xp(root, \emptyset, 0) \\
\mathcal{T}_{\mathcal{P}}^{\omega}(\mathbf{T})/_{v_k} &= fi\,xp(v, \mathcal{T}_{\mathcal{P}}^{\omega}(\mathbf{T})/_{v}, k).
\end{aligned}$$

*Proof of Lemma 8:* By induction on the tree (bottom-up):

- (Induction start) Let  $v_0$  be a leaf.  $\mathbf{T}_{v_0}$  is a tree with only a single node, so  $\mathbf{T}_{v_0}$  interpreted as a set of ground atoms only consists of unary atoms over  $v_0$ , and the same is true for  $X_0 \cup \mathbf{T}_{v_0}$ .  $fi\,xp(v_0, X_0, k = 0)$  ( $k$  must be 0 because there are no children.) initially sets  $X$  to  $X_0 \cup \mathbf{T}_{v_0}$  and then iterates  $X := \mathcal{T}_{\mathcal{P}}(X)$ . Note that since  $X$  at all times only contains unary atoms, only rules of types 1 and 4 can fire. When a fixpoint is reached (which must be reached eventually),  $X/v_0$  is returned. Thus,  $fi\,xp(v_0, X_0, 0) = \mathcal{T}_{\mathcal{P}}^{\omega}(X_0 \cup \mathbf{T}_{v_0})/_{v_0}$ , as claimed.
- (Induction step) Let  $v_0$  be a non-leaf node. We rewrite the pseudocode of Figure 5 by replacing the recursive calls to  $fi\,xp$  using the induction hypothesis. By our induction hypothesis, the pseudocode

```

X := X_0 ∪ T/{v_0, v_1, ..., v_n};
while no fixpoint reached do
  X := T_P(X) ∪
    T_P^ω(T_{v_1} ∪ X/v_1)/_{v_1} ∪ ... ∪
    T_P^ω(T_{v_n} ∪ X/v_n)/_{v_n};
return X/v_k;

```

obtained must equally compute  $fi\,xp(v_0, X_0, k)$ .

Clearly, the while-loop converges to a fixpoint, with  $\mathcal{T}_{\mathcal{P}}$  and set unions monotonically increasing and the set of ground atoms finite.

By the first part of Lemma 7, it is safe to decompose the computation of  $\mathcal{T}_{\mathcal{P}}^{\omega}(\mathbf{T}_{v_0} \cup X_0)/_{v_k}$  into the local processing of the region of  $\mathbf{T}$  consisting of  $v_0$  and its children  $v_1 \dots v_n$  and the separate processing of their subtrees  $\mathbf{T}_{v_1} \dots \mathbf{T}_{v_n}$  using the  $fi\,xp$  function.

It is essential to observe that  $X$  at all times only consists of the part of the tree structure that only involves nodes  $v_0, v_1, \dots, v_n$  as well as unary atoms over the same nodes. Because of the limited syntax of TMNF, rules of types 2 and 3 involving binary predicates can only (and only have to) unify with neighbouring nodes in the tree – either parents and their children or adjacent siblings. Therefore, all rules that can fire on  $\mathbf{T}_{v_0}$  involving  $v_0$  and its children will already fire on  $\mathbf{T}/_{v_0, v_1, \dots, v_n}$  given that the iterative computation of  $\mathcal{T}_{\mathcal{P}}^{\omega}(\mathbf{T}_{v_i} \cup X)/_{v_i}$  on the children ( $1 \leq i \leq n$ ) contributes the atoms over the child nodes  $v_i$  that can be computed in the subtrees  $\mathbf{T}_{v_i}$ , possibly given atoms over  $v_i$  computed locally on  $\mathbf{T}/_{v_0, v_1, \dots, v_n}$ .

Whenever  $\mathcal{T}_{\mathcal{P}}(X)$  adds a new ground atom over a child node  $v_i$  ( $i \geq 1$ ) to  $X$ ,  $\mathcal{T}_{\mathcal{P}}^{\omega}(\mathbf{T}_{v_i} \cup X/v_i)/_{v_i}$  is re-computed immediately afterwards to make sure that all consequences over node  $v_i$  from that atom that may be due to inferences in  $\mathbf{T}_{v_i}$  are added to  $X$ . Therefore, the while-loop only terminates when  $X = \mathcal{T}_{\mathcal{P}}^{\omega}(\mathbf{T}_{v_0} \cup X_0)/_{v_0, v_1, \dots, v_n}$ . Our claim follows.

Thus, Figure 5 indeed defines the function

$$\text{fixp}(v_0, X_0, k) = \mathcal{T}_{\mathcal{P}}^\omega(\mathbf{T}_{v_0} \cup X_0) / v_k.$$

Since  $\mathbf{T}_{root} = \mathbf{T}$ ,  $\mathcal{T}_{\mathcal{P}}^\omega(\mathbf{T}) / root = \text{fixp}(root, \emptyset, 0)$ . Moreover, by Lemma 7,  $\mathcal{T}_{\mathcal{P}}^\omega(\mathbf{T}) / v_k = \text{fixp}(v, \mathcal{T}_{\mathcal{P}}^\omega(\mathbf{T}) / v, k)$  (where  $v_k$  is the  $k$ -th child of  $v$ ).  $\square$

This provides us with an algorithm to compute  $\mathcal{T}_{\mathcal{P}}^\omega(\mathbf{T}) / v$  for an arbitrary  $v$  given its edge-path  $\Pi(v) = k_1 \cdots k_m$ . It simply starts from  $\text{fixp}(root, \emptyset, 0)$  and iteratively calls the  $\text{fixp}$  function for each of the segments  $k_i$  of  $\Pi(v)$ , always using the result of the previous call as the second and  $k_i$  as the third argument.

Note that differently from Figure 5, the second argument as well as the return value in the implementation of Figure 6 is of type “set of predicates” rather than “set of atoms”. Moreover, we make the recursive calls to  $\text{fixp}$  more intelligently, that is, only when new atoms have been derived at a child node and there is a chance of  $\text{fixp}$  inferring something new.

**Algorithm 4.1** (*Evaluation of monadic datalog*)

*Input:* A program  $\mathcal{P}$ , an instance  $\mathbf{I}$ , and a node  $v \in \Pi(V^I)$  (that is,  $v$  is a node of  $T(\mathbf{I})$  expressed as an edge-path).

*Output:* The set of IDB predicates for  $v$  that can be derived using  $\mathcal{P}$  on  $T(\mathbf{I})$ .

*Method:*

Let  $k_1 \dots k_m$  be the edge-path defining  $v$  ( $m \geq 0$ );

$w := root^I$ ;  $F := \text{fixp}(w, \emptyset, 0)$ ;  $i := 1$ ;

**while**  $i \leq m$  **do**

**begin**

$F := \text{fixp}(w, F, k_i)$ .

$w :=$  the  $k_i$ -th child of  $w$ ;

$i := i + 1$ ;

**end**;

Output  $F$ .

where the auxiliary function

$$\text{fixp} : V^I \times 2^{IDB(\mathcal{P})} \times \{0, 1, 2, 3, \dots\} \rightarrow 2^{IDB(\mathcal{P})}$$

is given in Figure 6.  $\square$

**Theorem 9.** *Given a TMNF program  $\mathcal{P}$ , an instance  $\mathbf{I}$ , and a node  $v \in T(\mathbf{I})$ , the set  $\{P(v) \in \mathcal{T}_{\mathcal{P}}^\omega(T(\mathbf{I}))\}$  can be computed in space  $O(\|\mathbf{I}\| * |IDB(\mathcal{P})|)$ .*

*Proof of Theorem 9:* The space requirements of the above-described algorithm are dominated by those of the  $\text{fixp}$  function. This function in turn only considers very localized regions of the instance at any point in time, namely a node  $v$  plus its children. It is recursive, and therefore  $\text{depth}(T(\mathbf{I}))$   $\text{fixp}$  activation records may have to be kept on the stack. The total space required may at most amount to the number

**function**  $\text{fixp}(\text{node } v, \text{ set of predicates } F, \text{ integer } k \geq 0)$   
**returns** set of predicates

**begin**

$F_0 := F \cup \{P \mid P(x) \leftarrow U(x). \text{ in } \mathcal{P}, U \text{ holds at } v\}$ ;

**let**  $n$  be the number of children of  $v$ ;

**for each** child  $v_i$  of  $v$  ( $1 \leq i \leq n$ ) **do**

$F_i := \text{fixp}(v_i, \emptyset, 0)$ ;

**while** no fixpoint of  $F_0$  reached **do**

**begin**

**if**  $P(x) \leftarrow P_1(x) \wedge P_2(x). \text{ in } \mathcal{P}$

**and**  $P_1, P_2 \in F_0$  **then**

$F_0 := F_0 \cup \{P\}$ ;

**if**  $P(x) \leftarrow P_0(x_0) \wedge \text{First-Child}(x_0, x). \text{ in } \mathcal{P}$ ,

$n \geq 1$  **and**  $P_0 \in F_0$  **then**

$F_1 := \text{fixp}(v_1, F_0 \cup \{P\}, 0)$ ;

**if**  $P(x) \leftarrow P_0(x_0) \wedge \text{First-Child}(x, x_0). \text{ in } \mathcal{P}$ ,

$n \geq 1$  **and**  $P_0 \in F_1$  **then**

$F_0 := F_0 \cup \{P\}$ ;

**if**  $P(x) \leftarrow P_0(x_0) \wedge \text{Next-Sibling}(x_0, x). \text{ in } \mathcal{P}$ ,

$P_0 \in F_i$  **and**  $1 \leq i < n$  **then**

$F_{i+1} := \text{fixp}(v_{i+1}, F_i \cup \{P\}, 0)$ ;

**if**  $P(x) \leftarrow P_0(x_0) \wedge \text{Next-Sibling}(x, x_0). \text{ in } \mathcal{P}$ ,

$P_0 \in F_i$  **and**  $1 < i \leq n$  **then**

$F_{i-1} := \text{fixp}(v_{i-1}, F_{i-1} \cup \{P\}, 0)$ ;

**end**;

**return**  $F_k$ ;

**end.**

Figure 6. A more detailed implementation of  $\text{fixp}$ .

of nodes of a path from the root to a leaf plus all of their siblings (times space to store a set of predicates of size at most  $|IDB(\mathcal{P})|$ ). Such a fragment of an instance has the same size no matter whether compressed or not. Thus, we obtain the space bound indicated in Theorem 9. Note that our algorithm proceeds in exactly the same way no matter whether the instance is a tree or compressed.  $\square$

Evaluating Core XPath (and thus TMNF) on compressed instances is also PSPACE-complete:

**Theorem 10.** *The evaluation problem for both monadic datalog and Core XPath over compressed instances is PSPACE-complete. Positive Core XPath over compressed instances is NP-complete.*

*Proof of Theorem 10:* PSPACE-membership for TMNF is

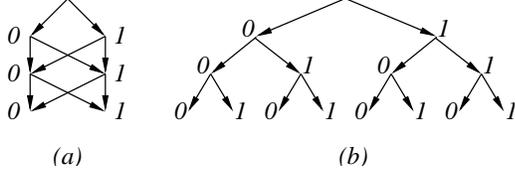


Figure 7. Instance for QSAT encoding (three variables).

asserted in Theorem 9, and Core XPath inherits this upper bound by Proposition 6. We show the PSPACE-hardness of both languages by providing a logspace-reduction from QSAT to the decision problem for Core XPath query evaluation.

Given a (closed) quantified boolean formula

$$Q_1 x_1 \cdots Q_n x_n \varphi$$

(where  $Q_1, \dots, Q_n \in \{\forall, \exists\}$  and  $\varphi$  is quantifier-free), the instance is as shown in Figure 7 for  $n = 3$  variables and has schema  $\sigma = \{0, 1\}$ .

The query is  $/\text{self}::*[Q_1 \text{ child}[\dots[Q_n \text{ child}[\varphi']]\dots]]$ , where  $\exists \text{ child}[\psi]$  rewrites into  $\text{child}::*[\psi]$  and  $\forall \text{ child}[\psi]$  rewrites into  $\text{not}(\text{child}::*[\text{not}(\psi)])$ .  $\varphi'$  is obtained from the boolean formula  $\varphi$  by substituting  $\wedge$ ,  $\vee$ , and  $\neg$  by “and”, “or”, and “not()”, respectively, and each variable  $x_i$  by  $\text{parent}::*/\dots/\text{parent}::*/\text{self}::1$ , where the number of applications of the parent axis is  $n - i$ .

It is not difficult to see that this is indeed a correct encoding of QSAT, i.e. the query selects the root node if and only if the QSAT formula is true.

For example, the query for the QSAT formula

$$\forall x_1 \exists x_2 (\neg x_1 \vee x_2) \wedge (x_1 \vee \neg x_2)$$

(which is true) is

$$/\text{self}::*[\text{not}(\text{child}::*[\text{not}(\text{child}::*[\varphi'])])]$$

with

$$\begin{aligned} \varphi' &= (\text{not}(\text{parent}::*/\text{self}::1) \text{ or } \text{self}::1) \text{ and} \\ &(\text{parent}::*/\text{self}::1 \text{ or } \text{not}(\text{self}::1)). \end{aligned}$$

The instance serves as a compressed version of a proof tree for the QSAT formula, where a node label (either 0 for “false” or 1 for “true”) at a node at depth  $i + 1$  in the tree denotes a valuation of propositional variable  $x_i$ . In our example, which can also be written as  $\forall x_1 \exists x_2 (x_1 \Leftrightarrow x_2)$ , the query selects the root node because both the left child of the root node labeled 0 and the right child labeled 1 ( $\forall x_1$ ) have a child ( $\exists x_2$ ) with their respective label ( $x_1 \Leftrightarrow x_2$ ).

The NP-hardness of positive Core XPath follows immediately from the encoding of the previous proof, in which

we can still encode SAT (i.e., all quantifiers are existential) without negation. Negation inside  $\varphi$  can be pushed down to the variables using De Morgan’s laws, where we can rewrite e.g.  $\text{not}(\text{parent}::*/\text{self}::1)$  as  $\text{parent}::*/\text{self}::0$ . For instance, we can write  $\varphi'$  equivalently as

$$(\text{parent}::*/\text{self}::0 \text{ or } \text{self}::1) \text{ and } (\text{parent}::*/\text{self}::1 \text{ or } \text{self}::0).$$

Membership in NP follows from the following observation. Let  $Q$  be a query and  $v$  be the node for which membership in the query result is to be checked. If negation is not present in  $Q$ , it can be matched in the instance  $I$  iff it can be matched on a small subtree of  $T(I)$  containing  $v$ . This “witness” tree is of polynomial size: it is at most as deep as  $I$  and is defined by a subset of  $\Pi(V^I)$  of cardinality at most the number of nodes of the query tree of  $Q$  (basically corresponding to the number of operations such as axis applications and node tests performed in  $Q$ ).

Indeed, no more than a subtree of this size can matter in the evaluation of such a query. The remaining parts of tree  $T(I)$  must be ignored and “jumped over” by the axis applications.

All we now need to do is guess such a witness tree and execute  $Q$  on it (for which we have a polynomial-time algorithm [17]).  $\square$

## 5. Binary Structures

The concise representation of multiple consecutive edges (just storing one edge together with a multiplicity) causes a number of problems. One way to get around these is to transform arbitrary instances into binary instances first and then only to work with these binary instances. In a binary instance we can store all edges explicitly, so there is no need for the multiplicities. For each instance  $I$  we define a binary instance  $B(I)$  as follows: We first replace a vertex with  $i$  children by an almost complete binary tree of height  $2^{\lceil \log(i) \rceil}$ , as indicated in Figure 8. The binary instance  $B(I)$  has an additional unary relation that contains all vertices of the original instance  $I$ . All other unary relations of  $I$  can be directly transferred to  $B(I)$ .

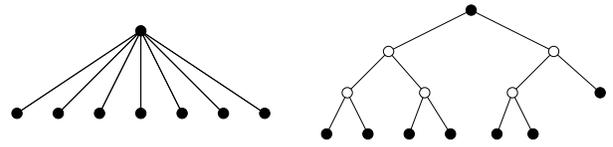


Figure 8.

We omit a formal definition of  $B(I)$ . The following proposition is crucial. Its proof is straightforward; Figure 9 illustrates why there will be a blow-up in size logarithmic in the maximum edge multiplicity.

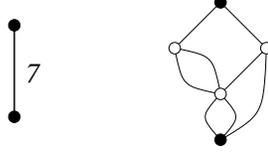


Figure 9.

**Proposition 11.** *There is an algorithm that, given an instance  $\mathbf{I}$ , computes a binary instance  $\mathbf{B}$  equivalent to  $\mathbf{B}(\mathbf{I})$  in time  $O(|\mathbf{I}| \cdot \log(\text{mult}(\mathbf{I})))$ .*

*Proof of Proposition 11:* Let  $\mathbf{I}$  be an instance and  $v \in V^{\mathbf{I}}$  with

$$\gamma(v) = w_1^{m_1} w_2^{m_2} \dots w_k^{m_k},$$

where  $w_1, \dots, w_k \in V^{\mathbf{I}}$  such that  $w_i \neq w_{i+1}$  for  $1 \leq i \leq k$ . Then we have  $m_i \leq \text{mult}(\mathbf{I})$  for  $1 \leq i \leq k$ . Let  $m = \sum_{i=1}^k k_i$  and  $h = \lceil \log(m) \rceil$ .

In the transformation from  $\mathbf{T}(\mathbf{I})$  to  $\mathbf{B}(\mathbf{T}(\mathbf{I}))$ , vertex  $v$  and its  $m$  children are replaced by a subtree  $\mathbf{S}$  of height  $h$  with  $m$  leaves. We claim that in the compressed binary instance  $\mathbf{M}(\mathbf{B}(\mathbf{I})) = \mathbf{M}(\mathbf{B}(\mathbf{T}(\mathbf{I})))$ , this subtree  $\mathbf{S}$  is compressed to a subinstance  $\mathbf{S}'$  with  $O(k \cdot \log(\text{mult}(\mathbf{I})))$  vertices.

To prove this claim, we assume without loss of generality that  $m = 2^h$ . If this is not the case, we can simply add a dummy vertex  $w_{k+1}$  and let  $\gamma(v) = w_1^{m_1} \dots w_k^{m_k} w_{k+1}^{2^h - m}$ . Then  $\mathbf{S}$  is a complete binary tree of height  $h$ . The  $i$ th level of  $\mathbf{S}$  and  $\mathbf{S}'$  consists of all vertices whose distance from  $v$  (the root of  $\mathbf{S}$ ) is  $i$ .

We observe the following:

- (1) On each level of  $\mathbf{S}$  there are at most  $2k - 1$  bisimilarity classes of vertices (i.e., at most  $2k - 1$  vertices that are pairwise not bisimilar). Thus each level of  $\mathbf{S}'$  contains at most  $2k - 1$  vertices.
- (2) For  $0 \leq \ell \leq h$ , the first  $\ell$  levels of  $\mathbf{S}$  contain  $2^{\ell+1} - 1$  vertices altogether. Thus the first  $\ell$  levels of  $\mathbf{S}'$  contain at most  $2^{\ell+1} - 1$  vertices.

Now let  $\ell = \lceil \log(k) \rceil$ . Then  $\mathbf{S}'$  contains at most

$$\begin{aligned} & (2^{\ell+1} - 1) + (h - \ell)(2k - 1) \\ & \leq (8k - 1) + (\log(m) - \log(k))(2k - 1) \\ & = (8k - 1) + \log\left(\frac{m}{k}\right)(2k - 1) \\ & = (8k - 1) + (\log(\text{mult}(\mathbf{I}))(2k - 1)) \\ & = O(k \cdot \log(\text{mult}(\mathbf{I}))). \end{aligned}$$

This completes the proof of the claim.

Noting that, given  $v$  and  $\gamma(v)$ , it is easy to compute  $\mathbf{S}'$  in time linear in the size of  $\mathbf{S}'$ , the statement of the proposition follows.  $\square$

**Remark 12.** Note that if we used a logarithmic cost model, then the logarithmic factor in Proposition 11 could be avoided, because storing  $m$  parallel edges would then require space and time  $\Theta(\log(m))$ .

While the binary instance computed by the algorithm in Proposition 11 may be larger by a factor of  $\log(\text{mult}(\mathbf{I}))$  than the original instance, it is not clear that it will be larger in practice. Indeed, after being minimised it may be exponentially smaller than the original instance, even if that was also minimal. The following example illustrates this.

**Example 13.** Let  $\sigma = \{P, Q\}$ . For  $\ell \geq 1$ , let  $\mathbf{I}_\ell$  be the following instance:  $\mathbf{I}$  has 3 vertices  $r, p, q$ .  $r$  is the root and  $p, q$  are leaves.  $\gamma$  is defined by  $\gamma(r) = (pq)^{2^\ell}$ ,  $\gamma(p) = \gamma(q) = \emptyset$ . Furthermore  $P = \{p\}$  and  $Q = \{q\}$ . The right hand side of Figure 10 shows  $\mathbf{I}_3$ .

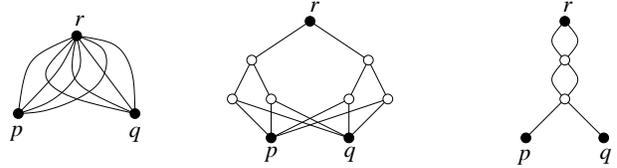


Figure 10. The instance  $\mathbf{I}_3$  of Example 13, the binary instance  $\mathbf{B}(\mathbf{I}_3)$ , and its minimisation  $\mathbf{M}(\mathbf{B}(\mathbf{I}_3))$ .

It is easy to see that  $\mathbf{I}_\ell$  is minimal and that  $\|\mathbf{I}_\ell\| \in O(2^\ell)$ . However, as Figure 10 illustrates, the minimal instance  $\mathbf{M}(\mathbf{B}(\mathbf{I}_\ell))$  bisimilar to the binary instance  $\mathbf{B}(\mathbf{I}_\ell)$  has size  $O(\ell)$ .

Translating queries in TMNF or MSO into queries over the corresponding binary instances is easy, so it may be well worth working with binary instances only.

Let  $\mathbf{B}^{-1}$  be the transformation that intuitively reverses the translation from unranked into binary instances and identifies each (non-auxiliary) node of  $\mathbf{B}(\mathbf{T})$  with the node of  $\mathbf{T}$  that it originates from. We assume that the schema of  $\mathbf{B}(\mathbf{T})$  contains a predicate  $V^{\mathbf{T}}$  that selects the nodes corresponding to the original nodes of  $\mathbf{T}$ , but not the auxiliary nodes introduced to obtain a binary tree.

**Lemma 14.** *Let  $\mathcal{P}$  be a TMNF program. Then there is a TMNF program  $\mathcal{P}'$  such that for all tree-instances  $\mathbf{T}$ ,*

$$\mathcal{T}_{\mathcal{P}}^\omega(\mathbf{T}) = \mathbf{B}^{-1}(\{P(v) \in \mathcal{T}_{\mathcal{P}'}^\omega(\mathbf{B}(\mathbf{T})) \mid P \in \text{IDB}(\mathcal{P})\}).$$

$\mathcal{P}'$  can be computed simultaneously in logarithmic space and linear time.

*Proof of Lemma 14:* We have to translate all occurrences of the predicates encoding the tree structure (i.e., the unary predicates *Root*, *Leaf* and *Last-Sibling* and the binary predicates *First-Child* and *Next-Sibling*) to obtain the correct behaviour over binary instances.

The *Root* and *Leaf* predicates remain the same on binary instances. The binary-instance version of *Last-Sibling* (called *Bin-Last-Sibling* below) can be defined by the following fixed program:

$$\begin{aligned}
Aux_0(x) &\leftarrow V^T(x). \\
Aux_1(x) &\leftarrow Aux_0(x_0) \wedge \\
&\quad Second-Child(x_0, x). \\
Aux_2(x) &\leftarrow Aux_0(x_0) \wedge \\
&\quad Last-Sibling(x_0). \\
Aux_1(x) &\leftarrow Aux_2(x_0) \wedge \\
&\quad First-Child(x_0, x). \\
Aux_0(x) &\leftarrow Aux_1(x) \wedge \overline{V^T}(x). \\
Bin-Last-Sibling(x) &\leftarrow Aux_1(x) \wedge V^T(x). \\
Bin-Last-Sibling(x) &\leftarrow Root(x).
\end{aligned}$$

The idea of this encoding is to start at a parent node and to find its last child by following the rightmost edge-path down, moving to the second child whenever one exists and only otherwise moving to the first child.

As all IDB predicates have to be unary and we thus cannot define a binary relation as a predicate, we have to rewrite rules containing binary predicates. There are four cases; we discuss two (the “forward” cases using TMNF rules of the second kind); the remaining two can be encoded similarly.

Consider the rule

$$P(x) \leftarrow P_0(x_0) \wedge First-Child(x_0, x).$$

In the binary version we essentially have to say that  $x$  is the first *First-Child*-descendent of  $x_0$  such that  $V^T(x)$  holds. Therefore, we replace this rule by the following. Note that for every rule we have to introduce a new predicate  $Aux-P_0$ .

$$\begin{aligned}
Aux-P_0(x) &\leftarrow P_0(x_0) \wedge First-Child(x_0, x). \\
Aux-P_0(x) &\leftarrow Aux-P_0(x_0) \wedge \overline{V^T}(x_0) \wedge \\
&\quad First-Child(x_0, x). \\
P(x) &\leftarrow Aux-P_0(x) \wedge V^T(x).
\end{aligned}$$

Accordingly, we proceed for a rule

$$P(x) \leftarrow P_0(x_0) \wedge Next-Sibling(x_0, x).$$

where we essentially check if  $x_0$  is a *Last-Sibling*, and if so, we move to its parent if that parent is not in  $V^T$ . Once we have found a non-*Last-Sibling* node, we go to its next sibling, and from there follow the *First-Child* path down. We select the first node in  $V^T$  we encounter. In particular, we replace the above rule by the following program fragment:

$$\begin{aligned}
Aux_1-P_0(x) &\leftarrow P_0(x_0). \\
Aux_2-P_0(x) &\leftarrow Aux_1-P_0(x) \wedge Last-Sibling(x).
\end{aligned}$$

$$\begin{aligned}
Aux_3-P_0(x) &\leftarrow Aux_2-P_0(x_0) \wedge First-Child(x, x_0). \\
Aux_3-P_0(x) &\leftarrow Aux_2-P_0(x_0) \wedge Second-Child(x, x_0). \\
Aux_1-P_0(x) &\leftarrow Aux_3-P_0(x) \wedge \overline{V^T}(x). \\
Aux_4-P_0(x) &\leftarrow Aux_1-P_0(x_0) \wedge First-Child(x, x_0). \\
Aux_5-P_0(x) &\leftarrow Aux_4-P_0(x_0) \wedge Second-Child(x_0, x). \\
Aux_5-P_0(x) &\leftarrow Aux_5-P_0(x_0) \wedge \overline{V^T}(x_0) \wedge \\
&\quad First-Child(x_0, x). \\
P(x) &\leftarrow Aux_5-P_0(x) \wedge V^T(x).
\end{aligned}$$

It is easy to see that this transformation can be done in logarithmic space and linear time.  $\square$

**Proposition 15.** *Let  $\mathbf{I}$  be an instance and  $Q$  a query in TMNF or MSO. Then, a query  $Q'$  in the same language can be computed in logarithmic space and linear time s.t.  $Q(\mathbf{I}) = Q'(\mathbf{B}(\mathbf{I}))$ .*

*Proof of Proposition 15:* For TMNF, the result follows immediately from Lemma 14, which makes an even stronger statement (namely that the translation correctly encodes all IDB predicates of  $\mathcal{P}$ , not just a single query predicate).

For MSO, we can directly take the encoding of a binary predicate  $B$  in TMNF and define a new binary relation  $\varphi_B$  in MSO with two free variables. In order to do that, we simply take the program  $\mathcal{P}_r$  encoding rule  $r$ :

$$P(x) \leftarrow P_0(x_0) \wedge B(x_0, x).$$

and interpret  $\leftarrow$  and  $\wedge$  as implication and logical conjunction in MSO. Let  $\bar{X} := IDB(\mathcal{P}_r)$  and

$$\begin{aligned}
\Psi_B(\bar{X}) &:= \bigwedge \{ \forall \bar{x} : r \mid r \in \mathcal{P}_{P_0, B, P}, \\
&\quad \bar{x} \text{ consists of the variables in } r \}
\end{aligned}$$

Now,

$$\varphi_B(x, y) := \forall \bar{X} ((x \in P_0 \wedge \Psi_B(\bar{X})) \rightarrow y \in P).$$

Given the TMNF program  $\mathcal{P}_{Last-Sibling}$  encoding *Last-Sibling* through the IDB predicate  $P$ , we proceed as above to obtain the FO-formula  $\Psi_{Last-Sibling}$ . Let  $\bar{X} := IDB(\mathcal{P}_{Last-Sibling})$ . We define

$$\varphi_{Last-Sibling}(x) := \forall \bar{X} (\Psi_{Last-Sibling}(\bar{X}) \rightarrow x \in P).$$

By replacing all occurrences of *First-Child*, *Next-Sibling* and *Last-Sibling* using the unary and binary relations defined in this way, we obtain the desired translated MSO query.  $\square$

Core XPath is not expressive enough to accommodate such a translation (consider, for example, the very simple query `/descendant::A/child::B`), but this is not a problem

as we translate all Core XPath queries into TMNF in our framework anyway.

### Experimental Results on Binary Instances.

We have carried out a number of experiments on XML corpora, comparing the compression achieved using our notion of binary instances with instances compressed using multiple edges.

To assess the merits of our framework for binary structures, we have extended the XML compression module of the XPath query engine discussed in [6] to create binary structures in the way discussed in Section 5. We have run this compressor on a number of standard corpora (see [6]). The benchmark results can be found in Figure 11. This table is interpreted as follows.

- The first (left-most) column provides the name and size of the corpus.
- The second column states the number of nodes in the XML tree. (Subtract one to obtain the number of edges.)
- The third and fourth columns provide the number of nodes and edges in the instance compressed using just bisimilarity, respectively.

For each corpus, we present in two rows

- in the upper row, the size of the compressed instance when XML node tags were ignored (i.e., only the bare tree structure is compressed, without any labeling information) and
- in the lower row, the size of the compressed instance with all XML node tags inserted.
- Columns 5 and 6 report the size of the instance (nodes and edges, respectively) compressed using edge multiplicities.
- Finally, columns 7 and 8 state the size of the instance (nodes and edges, respectively) compressed as a binary instance.

Our experiences with the prototype implementation of [6] are that in order to achieve good query performance with unranked instances, in the data structures we come up with, the amount of memory required for a node is approximately the same as for an edge. Thus, to estimate the size of an unranked instance, one best adds up the node and the edge count of the instance.

We observe that binary instances compress quite well; the compression seems to compare the more favorably the larger the corpora get. Together with the fact that binary instances can be represented more efficiently in memory (the data structures required to represent nodes and edges

have one degree of freedom less than unranked instances – each node only needs two fixed-size pointers to children, rather than an associated adjacency list), binary instances seem to be an interesting alternative to unranked instances compressed using multiple edges.

It is convenient to represent binary instances as relational structures in a slightly different way than arbitrary instances. We represent the edges of the DAG by two binary relations *First-Child* and *Second-Child*, representing the first-child and second-child relation.

In the following, we will give versions of our results for both unranked instances with edge multiplicities and binary instances.

## 6. Yet Another Query Automaton

The idea of querying XML-documents by tree automata is not new (e.g., [26, 19]). Some care needs to be taken, because in this context we are facing two problems usually not considered in classical language theory: Tree instances are not necessarily binary, but may be unranked, and queries select subsets of a tree and do not just accept or reject. Of course both of these problems can easily be resolved (and have been resolved in various ways, see e.g. [3, 12, 26]). We handle unranked trees simply by viewing them as binary trees under the *First-Child* and *Next-Sibling* relation. Our way of handling unary queries is similar to the approach proposed by Neven and Schwentick [26] in their *query automata*: certain states are *selecting states* that select the vertices in the answer of the query. However, beyond this superficial similarity the two querying mechanisms are quite different.

A non-deterministic (bottom-up) tree automaton is a tuple  $\mathfrak{A} = (Q, \Sigma, F, \delta)$ , where  $Q$  is the *state space*,  $\Sigma$  is the *alphabet*,  $F \subseteq Q$  is the set of *accepting states*, and

$$\delta : \Sigma \cup (Q \times \Sigma) \cup (Q \times Q \times \Sigma) \rightarrow 2^Q$$

the *transition function*. Tree automata work on binary trees in the usual way. Note that the transition function is defined in such a way that it accommodates leaves, nodes with one child, and nodes with two children. However, if a node has only one child, there is no way to distinguish between this child being a left or right child.

A  $\sigma$ -tree automaton is a nondeterministic tree-automaton  $\mathfrak{A} = (Q, \Sigma, F, \delta)$  with  $\Sigma = 2^\sigma$ , that is, a tree-automaton that runs on tree instances of schema  $\sigma$ . If  $T$  is a  $\sigma$ -tree instance and  $t \in V^T$ , we often write  $\Sigma(t)$  to denote the “symbol”  $\{R \in \sigma \mid t \in R^T\} \in \Sigma = 2^\sigma$ .

We can let the same tree automaton run on binary trees and on unranked trees, viewing the latter as binary trees with respect to the first-child and next-sibling relation. In

	$ V^T $	bisimilarity only		bisimilarity & edge multiplicities		bisimilarity on binary instances		include XML tags?
		$ V^{\mathcal{M}(T)} $	$ E^{\mathcal{M}(T)} $	$ V^{\mathcal{M}(T)} $	$ E^{\mathcal{M}(T)} $	$ V^{\mathcal{M}(T)} $	$ E^{\mathcal{M}(T)} $	
SwissProt (457.4 MB)	10,903,569	83,427	1,731,391	83,427	792,620	352,185	683,153	no
		85,712	1,751,930	85,712	1,100,648	371,650	721,429	yes
DBLP (103.6 MB)	2,611,932	321	272,573	321	171,820	72,020	143,951	no
		4,481	379,524	4,481	222,755	114,000	227,470	yes
TreeBank (55.8 MB)	2,447,728	323,256	909,875	323,256	853,242	504,759	966,591	no
		475,366	1,315,645	475,366	1,301,690	783,177	1,505,197	yes
OMIM (28.3 MB)	206,454	962	25,173	962	11,921	8,430	16,744	no
		975	25,173	975	14,416	8,549	16,949	yes
XMark (9.6 MB)	190,488	3,642	28,901	3,642	11,837	9,292	17,238	no
		6,692	39,180	6,692	27,438	15,080	27,329	yes
Shakespeare (7.9 MB)	179,691	1,121	40,855	1,121	29,006	15,806	31,381	no
		1,534	48,385	1,534	31,910	18,022	35,688	yes
Baseball (671.9 KB)	28,307	26	665	26	76	104	199	no
		83	1,378	83	727	469	834	yes
TPC-D (287.9 KB)	11,765	15	357	15	161	196	378	no
		53	361	53	261	303	507	yes

Figure 11. Degree of compression of benchmarked corpora.

the following, we treat binary trees in some details and then briefly explain how the approach extends to unranked trees.

**6.1. Binary instances.** A run of a  $\sigma$ -tree automaton  $\mathfrak{A} = (Q, \Sigma, F, \delta)$  on a binary tree instance  $T$  of schema  $\sigma$  is a mapping  $\rho : V^T \rightarrow Q$  such that:

- For all leaves  $t \in V^T$  we have  $\rho(t) \in \delta(\Sigma(t))$ .
- For all nodes  $t \in V^T$  with one child  $t_1$  we have

$$\rho(t) \in \delta(\rho(t_1), \Sigma(t)).$$

- For all nodes  $t \in V^T$  with two children  $t_1, t_2$  we have

$$\rho(t) \in \delta(\rho(t_1), \rho(t_2), \Sigma(t)).$$

The run  $\rho$  is *accepting* if  $\rho(\text{root}^T) \in F$ . The automaton  $\mathfrak{A}$  *accepts*  $T$  if there is an accepting run for  $\mathfrak{A}$  on  $T$ .

To be able to define unary queries, we need to enhance tree automata by an additional mechanism for selecting nodes.

**Definition 16.** A *selecting  $\sigma$ -tree automaton* ( $\sigma$ -STA) is a tuple  $\mathfrak{A} = (Q, \Sigma, F, \delta, S)$ , where  $(Q, \Sigma, F, \delta)$  is a  $\sigma$ -tree automaton and  $S \subseteq Q$  a set of *selecting states*.

The unary query defined by a  $\sigma$ -STA  $\mathfrak{A}$  maps every  $\sigma$ -tree  $T$  to the set

$$\mathfrak{A}(T) = \{v \in V^T \mid \text{every accepting run of } \mathfrak{A} \text{ on } T \text{ is in a selecting state at vertex } v\}.$$

At first sight, this querying mechanisms may seem a bit artificial. However, it turns out that STAs have good algorithmic properties and provide a nice unifying framework for the languages considered here. Still, requiring *all* accepting runs to be in a selecting state seems somewhat arbitrary. We shall see below that we may equivalently require *at least one* accepting run to be in a selecting state (cf. Corollary 19). Of course it would even be better if we could simply use *deterministic* tree automata.<sup>3</sup> The following example shows that this would not be sufficient.

**Example 17.** Let EVEN-DEPTH be the query of schema  $\emptyset$  defined by

$$\text{EVEN-DEPTH}(T) = \{v \in V^T \mid \text{depth}(v) \text{ is even}\}$$

where the depth of a vertex in a tree is the length of the path from the root to this vertex.

EVEN-DEPTH is clearly definable by an STA. It is not definable by a deterministic STA, though. To see this, just note that because there is only one run of an STA on every tree, for a query defined by a deterministic STA it only depends on the subtree below a vertex whether the vertex belongs to the answer set or not.

<sup>3</sup>A tree automaton is *deterministic* if  $\delta(\bar{q}, a)$  is a one-element set for all  $(\bar{a}, \bar{q}) \in \Sigma \cup (Q \times \Sigma) \cup (Q \times Q \times \Sigma)$ . The usual powerset construction shows that for every (nondeterministic) tree automaton  $\mathfrak{A}$  there is a deterministic tree automaton  $\mathfrak{A}'$  equivalent to  $\mathfrak{A}$ , in the sense that  $\mathfrak{A}$  and  $\mathfrak{A}'$  accept the same trees.

Neven and Schwentick's [26] query automata are deterministic, but the price for this is that a run of a query automaton may go up and down the tree several times. As the following result implies, both types of automata have the same querying power.

**Theorem 18 (see also [23]).** *A query (on binary tree-instances) is definable in MSO if, and only if, it is definable by an STA.*

*Proof of Theorem 18:* For the forward direction (i.e., that for every unary MSO query there exists an equivalent STA), we use the well known fact that a class of binary trees is definable in MSO if and only if it is recognisable by a tree automaton [9, 28]. It implies that for every MSO-formula  $\psi(X)$  of vocabulary  $\sigma \cup \{\text{First-Child}, \text{Second-Child}\}$  for some schema  $\sigma$  there is a binary  $\sigma \cup \{X\}$ -tree automaton  $\mathfrak{A}_\psi$  such that for all binary  $\sigma$ -tree instances  $\mathbf{T}$  and for all subsets  $U \subseteq V^{\mathbf{T}}$  we have

$$(\mathbf{T}, U) \models \psi(X) \iff \mathfrak{A}_\psi \text{ accepts } (\mathbf{T}, U).$$

Here  $(\mathbf{T}, U)$  denotes the expansion of  $\mathbf{T}$  to the  $\sigma \cup \{X\}$ -instance in which  $X$  is interpreted by  $U$ . Now let  $\varphi(x)$  be an MSO-formula with one free variable  $x$ . Let  $\psi(X) = \forall x(\varphi(x) \rightarrow Xx)$ . Then

$$\varphi(\mathbf{T}) = \bigcap \{U \subseteq V^{\mathbf{T}} \mid (\mathbf{T}, U) \models \psi(X)\}.$$

Suppose  $\mathfrak{A}_\psi = (Q, 2^{\sigma \cup \{X\}}, F, \delta)$ . Let  $\mathfrak{A}$  be the selecting  $\sigma$ -tree automaton  $(Q \times \{0, 1\}, 2^\sigma, F \times \{0, 1\}, \delta', Q \times \{1\})$ , where  $\delta'$  is defined by

$$\begin{aligned} \delta'((q_1, \varepsilon_1), (q_2, \varepsilon_2), a) = \\ (\delta(q_1, q_2, a) \times \{0\}) \cup (\delta(q_1, q_2, a \cup \{X\}) \times \{1\}) \end{aligned}$$

for all  $q_1, q_2 \in Q$ ,  $\varepsilon_1, \varepsilon_2 \in \{0, 1\}$ ,  $a \in 2^\sigma$  (similarly for  $\delta(a)$  and  $\delta'(q, a)$ ). Then accepting runs of  $\mathfrak{A}$  on a tree  $\mathbf{T}$  correspond to accepting runs of  $\mathfrak{A}_\psi$  on expansions  $(\mathbf{T}, U)$  of  $\mathbf{T}$ . Therefore,

$$\mathfrak{A}(\mathbf{T}) = \bigcap \{U \subseteq V^{\mathbf{T}} \mid \mathfrak{A}_\psi \text{ accepts } (\mathbf{T}, U)\} = \varphi(\mathbf{T}).$$

For the converse direction, let  $\mathfrak{A} = (Q, 2^\sigma, F, \delta, S)$  be a  $\sigma$ -STA with  $Q = \{q_1, \dots, q_k\}$ . A run  $\rho$  of  $\mathfrak{A}$  on a binary  $\sigma$ -tree instance  $\mathbf{T}$  can be described by a tuple  $(U_1, \dots, U_k)$  of subsets of  $V^{\mathbf{T}}$ , where  $v \in U_i$  if  $\rho(v) = q_i$ . It is easy to define a first-order formula  $\psi(X_1, \dots, X_k)$  saying that  $(X_1, \dots, X_k)$  describes an accepting run. Then

$$\varphi(x) = \forall X_1 \dots \forall X_k (\psi(X_1, \dots, X_k) \rightarrow \bigvee_{q_i \in S} X_i x)$$

defines the same query as  $\mathfrak{A}$ .  $\square$

**Corollary 19 (see also [23]).** *For every STA  $\mathfrak{A}$  there exists an STA  $\mathfrak{A}'$  such that for all binary tree instances  $\mathbf{T}$ ,*

$$\mathfrak{A}(\mathbf{T}) = \{v \in V^{\mathbf{T}} \mid \text{there exists an accepting run of } \mathfrak{A}' \text{ on } \mathbf{T} \text{ that is in a selecting state at vertex } v\}.$$

*Proof of Corollary 19:* Theorem 18 implies that the negation of a query defined by an STA can also be defined by an STA. Let  $\mathfrak{B} = (Q, 2^\sigma, F, \delta, S)$  be an STA defining the negation of the query defined by  $\mathfrak{A}$ . Let  $\mathfrak{A}' = (Q, 2^\sigma, F, \delta, Q \setminus S)$ . It is easy to see that  $\mathfrak{A}'$  has the desired property.  $\square$

Since monadic datalog is contained in MSO, Theorem 18 also implies that for every monadic datalog program there is an STA defining the same query. In the following example, we give a direct construction of such an automaton for a given datalog program, which of course is more efficient than the one going through MSO.

**Example 20.** Let  $\sigma$  be a schema and  $X_1, \dots, X_\ell \notin \sigma$ . Let  $\mathcal{P}$  be a TMNF-program of schema  $\sigma$ . Let  $X_1$  be the goal predicate. We define a  $\sigma$ -STA  $\mathfrak{A} = (Q, 2^\sigma, F, \delta, S)$  as follows: We let  $Q = 2^{\{X_1, \dots, X_\ell\}}$ ,  $F = Q$ , and  $S = \{q \in Q \mid X_1 \in q\}$ . To define the transition function, we first define a propositional Horn formula  $\Psi$  in the variables  $\sigma \cup \{X_i, X_i^1, X_i^2 \mid 1 \leq i \leq \ell\}$  (we consider relation names as propositional variables here) with the following clauses:

- If  $X_i(x) \leftarrow R(x)$  is a rule of  $\mathcal{P}$  then  $X_i \leftarrow R$  is a clause of  $\Psi$ .
- If  $X_i(x) \leftarrow X_j(x) \wedge X_k(x)$  is a rule of  $\mathcal{P}$  then  $X_i \leftarrow X_j \wedge X_k$  is a clause of  $\Psi$ .
- If  $X_i(x) \leftarrow X_j(y) \wedge \text{First-Child}(x, y)$  is a rule of  $\mathcal{P}$  then  $X_i \leftarrow X_j^1$  is a clause of  $\Psi$ .  
If  $X_i(x) \leftarrow X_j(y) \wedge \text{Second-Child}(x, y)$  is a rule of  $\mathcal{P}$  then  $X_i \leftarrow X_j^2$  is a clause of  $\Psi$ .
- If  $X_i(x) \leftarrow X_j(y) \wedge \text{First-Child}(y, x)$  is a rule of  $\mathcal{P}$  then  $X_i^1 \leftarrow X_j$  is a clause of  $\Psi$ .  
If  $X_i(x) \leftarrow X_j(y) \wedge \text{Second-Child}(y, x)$  is a rule of  $\mathcal{P}$  then  $X_i^2 \leftarrow X_j$  is a clause of  $\Psi$ .

Now for all  $q_1, q_2 \in Q$  and  $\tau \in 2^\sigma$  we let  $\Psi(q_1, q_2, \tau)$  be the formula in the variables  $\{X_1, \dots, X_\ell\}$  obtained from  $\Psi$  by replacing each variable  $X_j^i$ , for  $i = 1, 2$  and  $1 \leq j \leq \ell$ , by TRUE if  $X_j \in q_i$  and by FALSE otherwise, and by replacing each  $R \in \sigma$  by TRUE if  $R \in \tau$  and by FALSE otherwise. Then we let  $\delta(q_1, q_2, \tau)$  be the set of all satisfying assignments of  $\Psi(q_1, q_2, \tau)$ . Here a satisfying assignment is identified with the set of variables it sets TRUE.

For  $q_1 \in Q$  and  $\tau \in 2^\sigma$  we define a formula  $\Psi(q_1, \tau)$  as  $\Psi(q_1, q_2, \tau)$  above, just replacing all variables  $X_j^2$  by FALSE. Then we let  $\delta(q_1, \tau)$  be the set of all satisfying assignments of  $\Psi(q_1, \tau)$ . For  $\tau \in 2^\sigma$  we define  $\delta(\tau)$  similarly.

This completes the definition of the automaton  $\mathfrak{A}$ . It is not hard to prove that  $\mathfrak{A}$  defines the same query as the program  $\mathcal{P}$ .

We now show how to evaluate STA-queries, that is, queries defined by STAs, first on tree-instances and then on compressed instances. The algorithms combine ideas from [12] with the partial decompression methods that we also use to evaluate XPath and monadic datalog queries.

**Proposition 21.** *The evaluation problem for STA-queries on binary tree instances can be solved in time  $O(s^3 \cdot n)$ , where  $s$  is the number of states of the input automaton and  $n$  the number of vertices of the input instance.*

*Proof of Proposition 21:* Let  $\mathfrak{A} = (Q, 2^\sigma, F, \delta, S)$  be an STA and  $\mathbf{T}$  a binary  $\sigma$ -tree instance.

For every  $t \in V^{\mathbf{T}}$ , let  $reach(t)$  be the set of all states  $q$  such that there is a run  $\rho$  of  $\mathfrak{A}$  on  $\mathbf{T}$  with  $\rho(t) = q$ . The states in  $reach$  are called the *reachable states* at  $t$ . Clearly, the mapping  $reach : V^{\mathbf{T}} \rightarrow 2^Q$  can be computed in time  $O(s^3 \cdot n)$  in a bottom-up pass of the tree.

For every  $t \in V^{\mathbf{T}}$ , let  $succ(t)$  be the set of all states  $q$  such that there is an *accepting* run  $\rho$  of  $\mathfrak{A}$  on  $\mathbf{T}$  with  $\rho(t) = q$ . The states in  $succ$  are called the *successful states* at  $t$ . Observe that

$$succ(root^{\mathbf{T}}) = reach(root^{\mathbf{T}}) \cap F,$$

and that for a siblings  $t_1, t_2 \in V^{\mathbf{T}}$  with parent  $t \in V^{\mathbf{T}}$ ,

$$succ(t_1) = \{q_1 \in reach(t_1) \mid \exists q_2 \in reach(t_2), q \in succ(t) : q \in \delta(q_1, q_2, \Sigma(t))\}.$$

$succ(t_2)$  can be obtained analogously. Similarly,  $succ(t_1)$  can be obtained from  $succ(t)$  for an only-child  $t_1$ . Using this, it is easy to see that given  $reach$ , the mapping  $succ : V^{\mathbf{T}} \rightarrow 2^Q$  can be computed in time  $O(s^3 \cdot n)$  in a top-down pass of the tree.

The vertices in  $\mathfrak{A}(\mathbf{T})$  are precisely the  $t \in V^{\mathbf{T}}$  for which  $succ(t) \subseteq S$ .  $\square$

**Remark 22.** Note that storing the transition relation of an automaton with  $s$  states requires space  $\Omega(s^2)$  and  $O(s^3)$ . So the factor  $s^3$  in the running time of the previous algorithm is not as bad as it looks. Indeed, a closer analysis shows that the running time can be improved to  $O(m \cdot n)$ , where  $m$  is the size of the encoding of the automaton.

In the remainder of this section, whenever we have an automaton  $\mathfrak{A}$  run on a (compressed) instance  $\mathbf{I}$ , the actual meaning is that  $\mathfrak{A}$  runs on the unfolded tree-instance  $\mathbf{T}(\mathbf{I})$ , even if our results will be based on techniques that usually do not require complete decompression to answer

queries. Thus, we continue to address the problem of evaluating queries on (compressed) trees, rather than studying automata for directed acyclic graphs.

**Theorem 23.** *The evaluation problem for STA-queries on binary instances can be solved in time  $O(2^{s+3\log s} \cdot n)$ , where  $s$  is the number of states of the input automaton and  $n$  the number of vertices of the input instance.*

*Proof of Theorem 23:* We proceed as on tree instances: In a first bottom up pass we compute the sets of reachable states at all vertices, and in a second, top-down pass we compute the successful states. The problem is that a reachable state  $q \in reach(v)$  may be successful on some path from  $v$  to the root and not successful on some other path. If this happens, we have to create a copy  $v'$  of vertex  $v$  and put  $q$  into  $succ(v')$ , but not into  $succ(v)$  (or vice versa). In the worst case, we have to create a copy of every vertex for every subset of the state space, which explains why the running time is exponential in  $s$ .

Let us make this precise. Let  $\mathfrak{A} = (Q, 2^\sigma, F, \delta, S)$  be a  $\sigma$ -STA and  $\mathbf{I}$  a binary  $\sigma$ -instance. Let  $\mathbf{T} = \mathbf{T}(\mathbf{I})$  be the tree instance bisimilar to  $\mathbf{I}$ . Recall that  $\Pi : V^{\mathbf{I}} \rightarrow 2^{V^{\mathbf{T}}}$  maps each vertex  $v \in V^{\mathbf{I}}$  to the set of vertices of  $\mathbf{T}$  it corresponds to.

For every  $t \in V^{\mathbf{T}}$  we define the reachable states  $reach(t)$  and successful states  $succ(t)$  at  $t$  as in the proof of Proposition 21. Since  $reach(t)$  only depends on the subtree below  $t$ , for all  $v \in V^{\mathbf{I}}, t, t' \in \Pi(v)$  we have  $reach(t) = reach(t')$ . Let  $reach(v) = reach(t)$ . Observe that the mapping  $reach : V^{\mathbf{I}} \rightarrow 2^Q$  can be computed in  $O(s^3 \cdot n)$  in a bottom-up pass of the instance  $\mathbf{I}$ . (Compression does not matter here.)

However, there may be  $v \in V^{\mathbf{I}}, t, t' \in \Pi(v)$  such that  $succ(t) \neq succ(t')$ . For  $v \in V^{\mathbf{I}}$  we let

$$\mathcal{S}(v) = \{succ(t) \mid t \in \Pi(v)\}.$$

In the top-down pass over the instance  $\mathbf{I}$  we compute a new instance  $\mathbf{J}$  and mappings  $f : V^{\mathbf{J}} \rightarrow V^{\mathbf{I}}$  and  $succ : V^{\mathbf{J}} \rightarrow 2^\sigma$  such that:

- (1)  $\mathbf{J}$  is bisimilar to  $\mathbf{I}$ ,
- (2) For all  $w \in V^{\mathbf{J}}$  we have  $w \approx f(w)$  and thus  $\Pi(w) = \Pi(f(w))$ .
- (3) For all  $w \in V^{\mathbf{J}}$  and  $t \in \Pi(w)$  we have  $succ(t) = succ(w)$ .
- (4) For every  $v \in V^{\mathbf{I}}$  and every  $X \in \mathcal{S}(v)$  there exists exactly one vertex  $w \in V^{\mathbf{J}}$  such that  $f(w) = v$  and  $succ(w) = X$ .

We start by creating a vertex  $root^{\mathbf{J}}$  and let  $succ(root^{\mathbf{J}}) = reach(root^{\mathbf{I}}) \cap F$  and  $f(root^{\mathbf{J}}) = root^{\mathbf{I}}$ .

Now suppose that  $w$  is a vertex of  $w$  that we have already created, and let  $v = f(w)$ , and  $X = \text{succ}(w)$ . Furthermore, suppose that  $v$  has children  $v_1$  and  $v_2$  and that we have not yet created any children of  $w$ . Let

$$X_1 = \{q_1 \in \text{reach}(v_1) \mid \exists q_2 \in \text{reach}(v_2), q \in X : q \in \delta(q_1, q_2, \Sigma(v))\}.$$

Then for all  $t \in \Pi(w)$  with children  $t_1$  and  $t_2$  we have  $X_1 = \text{succ}(t_1)$ . To see this, note that  $\text{reach}(t_1) = \text{reach}(v_1)$ ,  $\text{reach}(t_2) = \text{reach}(v_2)$ ,  $\text{succ}(t) = X$ , and recall the corresponding claim in the proof of Proposition 21.

If there already is a vertex  $w_1 \in V^J$  such that  $f(w_1) = v_1$  and  $\text{succ}(w_1) = X_1$ , we make  $w_1$  the first child of  $w$ . Otherwise, we create a new vertex  $w_1$ , let  $f(w_1) = v_1$  and  $\text{succ}(w_1) = X_1$ , and make  $w_1$  the first child of  $w$ . The second child of  $w$  is created analogously. If the vertex  $v = f(w)$  only has one child, we proceed similarly.

With an appropriate dictionary data structure that stores, for each vertex  $v \in V^I$ , the set of all pairs  $w$  with  $f(w) = v$  and allows, given  $X$ , to find the  $w$  with  $\text{succ}(w) = X$  (if it exists) in time  $O(s)$ , the children of a vertex of  $w$  can be created in time  $O(s^3)$ . Since  $J$  has at most  $2^s \cdot n$  vertices, the overall running time to compute  $J$  and  $\text{succ}$  is  $O(s^3 \cdot 2^s \cdot n) = O(2^{s+3\log s} \cdot n)$ .

The output of the algorithm is  $J$  together with the set of all vertices  $v \in V^J$  with  $\text{succ}(v) \subseteq S$ .  $\square$

**Remark 24.** It is not hard to see the decision version of the evaluation problem for STA-queries is in polynomial time, because for a given  $\sigma$ -STA we can easily define a  $\sigma \cup \{X\}$ -tree automaton  $\mathfrak{A}'$  such that for every  $\sigma$ -tree instance  $T$  and every vertex  $t \in V^T$  we have  $t \in \mathfrak{A}'(T)$  if and only if  $\mathfrak{A}'$  accepts  $(T, \{t\})$ . The latter can be tested in time  $O(s^3 \cdot n)$ , because an automaton accepts a tree if and only if there is an accepting reachable state at the root.

**Corollary 25.** *The evaluation problem for MSO-queries on binary instances parameterized by the size of the input formula is fix-parameter tractable.*

Let  $\exists\text{MSO}$  be the set of MSO-formulas of the form  $\exists X_1, \dots, X_k \varphi(X_1, \dots, X_k)$  where  $\varphi$  is first-order.

**Theorem 26.** *The evaluation problem for  $\exists\text{MSO}$ -queries on binary instances is NEXPTIME-complete. Evaluation of MSO queries can be done in EXPSPACE.*

*Proof of Theorem 26:* For the upper bounds recall that the tree-instance  $T(I)$  has size  $\leq 2^{|V^I|}$ . To decide whether  $\psi = \exists \bar{X} \varphi(\bar{X})$  holds, we guess  $k$  sets  $A_i$  and check whether  $T(I) \models \varphi(\bar{A})$ . This can be done in the straightforward way in time  $O(\|T(I)\|^{\|\varphi\|}) = O(2^{\|\varphi\| \cdot |V^I|})$ .

In case of full MSO, we use the standard PSPACE algorithm on  $T(I)$  [29], which, in total, yields EXPSPACE as upper bound.

For the hardness, assume that  $M$  is a fixed nondeterministic Turing machine with an NEXPTIME-complete acceptance problem (such a machine exists). We assume that  $M$  has alphabet  $\Delta = \{a_1, \dots, a_l\}$  and works on input  $w \in \Delta^n$  in time  $2^{n^k} - 1$ , for some  $k \geq 1$ . Let  $Q = \{q_1, \dots, q_m\}$  be the state space of  $M$ . Then configurations are words  $v \in \Delta^* Q \Delta^+$  of length  $2^{n^k}$  meaning that if  $v_i \in Q$  then  $M$  is in state  $q = v_i$  and currently reading the  $i$ th cell of its (input/work) tape (or equivalently,  $M$  reads  $v_{i+1}$ ).

From input  $w \in \Delta^n$  for  $M$  we construct a binary instance  $I$  with  $|V^I| = O(n^k)$  and an  $\exists\text{MSO}$ -formula  $\varphi$  with  $\|\varphi\| = n^{O(k)}$  such that  $M$  accepts  $w$  iff  $T(I) \models \varphi$ . Let  $I$  be the structure from Figure 7 with paths of length  $2 \cdot n^k + 1$  and let  $T$  denote its tree version  $T(I)$ . Note here that  $T$  has exactly  $2^{2 \cdot n^k} = 2^{n^k} \cdot 2^{n^k}$  leaves.

Recall that elements of  $T$  are edge-paths. We call every node  $t \in T$ , where  $t$  is a  $n^k$ -length edge-path, a *base*. There are  $2^{n^k}$  bases and each subtree rooted at a base has  $2^{n^k}$  leaves. Leaves and bases are each totally ordered by the lexicographic ordering of the respective edge-paths (the first digit is the most significant one). For convenience, we denote both of them by  $\prec$ . For  $0 \leq i \leq 2^{n^k} - 1$ , by  $t(i)$  we denote the  $i$ -th base of  $V^T$ . Analogously, let  $t^s(i)$  be the  $i$ -th leaf of the subtree rooted at base  $s \in V^T$ .

There is an FO-formula  $\psi(x, y)$  with  $\|\psi(x, y)\| = O(n^{2k})$  such that  $T \models \psi(t, t')$  iff  $t'$  is the direct successor of  $t$  (w.r.t.  $\prec$ ). For that, we simply guess two paths from the root to  $t$  and  $t'$  and check, if the edge-paths linking them are direct successors. Analogously, we can define a formula  $\psi(x_1, y_1, x_2, y_2)$  saying that  $x_1, x_2$  are bases,  $y_1$  and  $y_2$  are in  $T_{x_1}, T_{x_2}$  resp. and the path from  $x_2$  to  $y_1$  is the direct successor of the path from  $x_1$  to  $y_2$  (w.r.t.  $\prec$ ).

Recall that  $T_s$ , for a  $s \in V^T$ , denotes the subtree of  $T$  rooted at  $s$ . We encode a configuration  $v \in \Delta^* Q \Delta^+$  (of size  $2^{n^k}$ ) relative to base  $s$  as the  $(m + l)$ -tuple  $\bar{C} \subseteq V_s^T$  such that the following conditions hold:

- (1) if  $t \in C_\nu$  for some  $i$ , then  $t$  is a leaf
- (2) if  $v_i \in Q$  then  $t^s(i) \in C_\nu$  iff  $v_i = q_\nu$ , ( $\nu = 1, \dots, m$ )
- (3) if  $v_i \in \Delta$  then  $t^s(i) \in C_\nu$  iff  $v_i = a_{\nu-m}$ , ( $\nu = m + 1, \dots, m + l$ )

From now on we do not distinguish explicitly between configurations and their encodings. It is straightforward to define FO-formulas  $\varphi_{\text{start}}(\bar{X}, x)$  and  $\varphi_{\text{acc}}(\bar{X}, x)$ , which hold, if  $x$  is a base and  $\bar{X}$  relative to  $x$  is the *start* configuration or a *accepting* configuration, respectively. The intention behind these definition is to guess a tuple  $\bar{C}$  over the leaves of  $T$  and to check, whether the configurations relative

to the bases in  $\mathbf{T}$  form an accepting run. Here it is crucial to observe that  $\psi(x_1, y_1, x_2, y_2)$  allows us to identify consecutive positions  $(y_1, y_2)$  in two different configurations (relative to  $x_1, x_2$ ).

Next, we define  $\varphi_{next}(\bar{X}, x, y)$ , which essentially says that  $x$  and  $y$  are bases and direct  $\prec$ -successors, and that  $\bar{X}$  relative to  $x$  and  $\bar{Y}$  relative to  $y$  are configurations. We let denote these configurations by  $C_x$  and  $C_y$ . Now  $\varphi_{next}(\bar{X}, x, y)$  says that  $C_x$  “equals”  $C_y$ , or that we can pass from  $C_x$  to  $C_y$  in a single step. In particular, this means that there are  $i \leq n^k, \nu \leq m$  such that  $t^x(i) \in X_\nu$ , i.e. in configuration  $C_x$  we have the situation that  $M$  is state  $q_\nu$  and reads letter  $a_j$  ( $j$  such that  $t^x(i+1) \in X_{m+j}$ ). Since, using  $\psi$ , we are able to directly address successors of  $t^x(i)$  and  $t^y(i)$ , it is easy to devise a formula  $\varphi'(\bar{X}, x, y)$  saying that the neighborhoods around  $t^x(i)$  and  $t^y(i)$  conform to the transition relation of  $M$ .

The sought formula now looks as follows. Let

$$\varphi = \exists \bar{X} (\varphi_{start}(\bar{X}, 0) \wedge \varphi_{acc}(\bar{X}, 1) \wedge \forall \text{bases } x, y (\psi(x, y) \rightarrow \varphi_{next}(\bar{X}, x, y))),$$

where we use 0 and 1 to denote the first and last base w.r.t.  $\prec$ , respectively. It is easy to see that  $\varphi$  behaves as it is supposed to. More precisely, it states that there is a sequence  $v_1, \dots, v_{2^{n^k}}$  of configurations such that  $v_1, v_{2^{n^k}}$  is a starting, accepting configuration, respectively. Furthermore, it states that for each  $i < 2^{n^k}$  either  $v_{i+1}$  is equal to  $v_i$  or a direct successor of  $v_i$  (here w.r.t. the transition relation defined by  $M$ ). Together, this means that  $M$  on input  $w$  has an accepting run.

The size of  $\varphi$  is bounded by  $n^{O(k)}$ , caused by the subformulas needed to identify  $\prec$ -successors of leaves and bases. The remaining subformulas solely depend on the (fixed) machine  $M$ .  $\square$

If we convert a monadic datalog program into an STA in the way described in Example 20 and then use Theorem 23 to evaluate the query, the resulting algorithm is doubly exponential in the size of the datalog program, which is much worse than the direct algorithms we saw before. We shall now introduce a restricted version of our STA, which is somewhat closer to monadic datalog and admits more efficient query evaluation.

**Definition 27.** A *weak selecting  $\sigma$ -tree automaton* ( $\sigma$ -WSTA) is a tuple  $\mathfrak{A} = (Q, \Sigma, F, \delta, S, \preceq)$ , where  $(Q, \Sigma, F, \delta, S)$  is a  $\sigma$ -STA and  $\preceq$  a partial order on  $Q$  s.t.

- $F$  is downward closed w.r.t.  $\preceq$ .
- $S$  is upward closed w.r.t.  $\preceq$ .
- For all  $(\bar{q}, a) \in \Sigma \cup (Q \times \Sigma) \cup (Q \times Q \times \Sigma)$  and  $r, r' \in \delta(\bar{q}, a)$  there is an  $r'' \in \delta(\bar{q}, a)$  such that  $r'' \preceq r, r'$ .

- $\delta$  is monotone with respect to  $\preceq$  in the sense that for all  $(\bar{q}, a), (\bar{q}', a) \in (Q \times \Sigma) \cup (Q \times Q \times \Sigma)$ , if  $\bar{q} \preceq \bar{q}'$  (componentwise) then for every  $r' \in \delta(\bar{q}', a)$  there is an  $r \in \delta(\bar{q}, a)$  such that  $r \preceq r'$ .

**Example 28.** The automaton constructed from a monadic datalog program in Example 20 with set-inclusion as a partial order on the state space  $2^{\{X_1, \dots, X_n\}}$  is a WSTA. To see this, just note that the set of satisfying assignments of a propositional Horn formula is closed under intersection and monotone with respect to the Boolean constants in the formula.

**Proposition 29.** *For every STA there is a WSTA that defines the same query.*

*Proof of Proposition 29:* This follows from the fact that every MSO-query is definable by a monadic datalog program [16], Theorem 18, and the previous example.  $\square$

Note that if  $\mathfrak{A} = (Q, \Sigma, F, \delta, S, \preceq)$  is a  $\sigma$ -WSTA and  $\mathbf{T}$  a binary  $\sigma$ -tree instance, then  $\preceq$  induces a partial order on the runs of  $\mathfrak{A}$  on  $\mathbf{T}$ , which we also denote by  $\preceq$ , and which is defined by  $\rho \preceq \rho'$  if  $\rho(t) \preceq \rho'(t)$  for all  $t \in V^{\mathbf{T}}$ .

**Lemma 30.** *Let  $\mathfrak{A} = (Q, \Sigma, F, \delta, S, \preceq)$  be a  $\sigma$ -WSTA. Then for every binary  $\sigma$ -tree instance  $\mathbf{T}$  there is a unique minimal run  $\rho_{\min}$  of  $\mathfrak{A}$  on  $\mathbf{T}$  with respect to the order  $\preceq$  on the runs. Moreover, if  $\mathfrak{A}$  accepts  $\mathbf{T}$  then  $\rho_{\min}$  is an accepting run and*

$$\mathfrak{A}(\mathbf{T}) = \{t \in V^{\mathbf{T}} \mid \rho_{\min}(t) \in S\}.$$

*Proof:* In a straightforward bottom-up fashion we can construct for any two runs  $\rho, \rho'$  of  $\mathfrak{A}$  on  $\mathbf{T}$  a run  $\rho''$  such that  $\rho'' \preceq \rho, \rho'$ . Since there are only finitely many runs, this implies the existence of a unique minimal run.

If there is some accepting run  $\rho$ , then all runs  $\rho' \preceq \rho$  are also accepting, because  $F$  is downward closed. Thus in particular,  $\rho_{\min}$  is accepting.  $\mathfrak{A}(\mathbf{T}) = \{t \in V^{\mathbf{T}} \mid \rho_{\min}(t) \in S\}$  follows from the fact that  $S$  is upward closed w.r.t.  $\preceq$ , which implies that for runs  $\rho \preceq \rho'$ , every vertex selected by  $\rho$  is also selected by  $\rho'$ .  $\square$

**Theorem 31.** *The evaluation problem for WSTA-queries on binary instances can be solved in time  $O(s^3 \cdot n)$ , where  $s$  is the number of states of the input automaton and  $n$  the number of vertices of the input instance.*

*Proof of Theorem 31:* We essentially proceed as in the proof of Theorem 23 using Lemma 30 to restrict attention to success states instead of success sets. Let  $\mathfrak{A} = (Q, \Sigma, F, \delta, S, \preceq)$  be a  $\sigma$ -WSTA and  $\mathbf{I}$  a binary  $\sigma$ -instance. Let  $\mathbf{T} = \mathbf{T}(\mathbf{I})$  be the tree instance bisimilar to  $\mathbf{I}$ . For  $t \in V^{\mathbf{T}}$  define  $reach(t)$  and  $succ(t)$  as before. Furthermore, let  $succ'(t)$  be the minimal state of  $succ(t)$  (w.r.t.  $\preceq$ ) and let

$\mathcal{S}'(v) = \{succ'(t) \mid t \in \Pi(v)\}$ . We proceed as in the case of general STAs, but compute minimal states instead of sets of states. By the previous lemma this is sufficient to determine the sought result.

Again we start with a bottom-up pass computing  $reach(t)$  for all  $t \in V^I$ . Here note that  $reach(t) = reach(t')$  for all  $t, t' \in \Pi(v)$  for an arbitrary  $v \in V^I$ . Then we go top-down and gradually compute a new instance  $\mathbf{J}$  and mappings  $f : V^J \rightarrow V^I$  and  $succ' : V^J \rightarrow Q$  that satisfy the conditions (1) to (4) with  $succ$  and  $\mathcal{S}$  replaced by  $succ'$  and  $\mathcal{S}'$ , respectively. Note that, since  $F$  is downward closed,  $\rho(t) := succ'(v)$  for  $t \in \Pi(v)$  is the sought minimal accepting run  $\rho$  of  $\mathfrak{A}$  on  $\mathbf{T}$ .

We start by creating a vertex  $root^J$  and let  $succ'(root^J) = reach(root^I) \cap F$  and  $f(root^J) = root^I$ .

If  $w$  is a vertex of  $V^J$  that has already been created, let  $v = f(w)$  and  $q = succ'(w)$ . Suppose that  $v$  has children  $v_1, v_2$  and let  $X_1 = \{q_1 \in reach(v_1) \mid \exists q_2 \in reach(v_2) \text{ s.t. } q \in \delta(q_1, q_2, \Sigma(v))\}$ . Let  $q_1$  be the minimum of  $X_1$ .

Now we check if there is a  $w_1 \in V^J$  with  $f(w_1) = v_1$  and  $succ'(w_1) = q_1$ . If so, we make  $w_1$  the first child of  $w$ . Otherwise, we create a new  $w_1$  in  $\mathbf{J}$  and define  $f(w_1) = v_1$ ,  $succ'(w_1) = q_1$  and make this  $w_1$  the first child of  $w$ .

Now we see that for all  $t \in \Pi(w)$  with children  $t_1, t_2$  we have  $q_1 = succ'(t_1)$ . This shows that the procedure yields the intended result. Then we proceed analogously for the second child  $v_2$ .

But in contrast to the general case, now we create at most  $s$  new vertices for one vertex in  $\mathbf{I}$ , hence we have at most  $s \cdot n$  vertices in  $V^J$ .

An ad-hoc implementation yields an  $s^4$  factor. This can be improved as follows: at each node  $v \in V^I$  with children  $v_1, v_2$  we compute the sets  $pre_1(v, q) = \{q_1 \mid \exists q_2 \in reach(v_2), q \in \delta(q_1, q_2, \Sigma(v))\}$ . Note that this set coincides with  $X_1$ , computed for  $q = succ'(w)$  as defined above.

Analogously, we compute the sets  $pre_2(v, q)$ . Note that, using suitable data structures, this can be done in a single loop over all tuples  $(q_1, q_2, \Sigma(v), q) \in \delta$ , hence requires  $O(s^3)$  steps for each  $v \in V^I$ . Then we simply look up these tables to determine the values  $succ'(t)$  (we find the minimum in  $O(s)$  steps).

Together, this gives the  $O(s^3 \cdot n)$  time bound.  $\square$

**6.2. Unranked tree automata.** Even though we can avoid unranked instances entirely if we transform them to binary instances as described in Section 5, we think it is still worthwhile to briefly discuss STAs on unranked instances. We have mentioned earlier that we run a tree automaton on an unranked tree by running it on the corresponding binary tree induced by the *First-Child* and *Next-Sibling* relation. Note

that in this binary tree, the last siblings are precisely those that only have one (first) child.

To distinguish them clearly from the “binary” automata considered in the previous section, we call the automata considered here *unranked tree automata*. Let us emphasise, however, that the automata themselves are the same, they only run in a different way.

Once we have defined an appropriate notion of run and acceptance for unranked tree automata, we can define *unranked STAs* and *unranked WSTAs*. Then the basic results stated for binary automata in the previous section have analogous versions for unranked automata. In particular, every MSO-query on unranked trees (viewed as  $\sigma \cup \{Root, Leaf, Last-Sibling, First-Child, Next-Sibling\}$ -structures) is definable by an STA and vice versa. Moreover, monadic datalog programs over unranked trees can be translated into WSTAs whose state space is the set of all IDB predicates.

We only state the unranked version of the main result. The bounds on the running time are weaker, which is mainly due to the compact representation of multiple edges in compressed instances (cf. Section 2.2).

**Theorem 32.** *The evaluation problem for*

- (1) *STA-queries can be solved in time  $O(2^{3s} + 2^s \cdot m)$ .*
- (2) *WSTA-queries can be solved in time  $O(2^{3s} + s \cdot m)$ .*

*Here  $s$  is the number of states of the input automaton and  $m$  the size of the input instance.*

*Proof of Theorem 32:* Basically, we proceed as in the binary case. When computing the reachable and successful states, it is convenient to associate them with edges of the instance instead of vertices. For example,  $reach(e)$  for an edge  $e = (v, w)$  represents the set  $reach(\Pi(w))$ . The reason we have to do this is that vertices may play different roles as children of different parents — they may be the first child of some, but not for all parents. If we transfer the algorithms in a straightforward way, we obtain running times of  $O(2^{s+3\log s} \cdot |E|)$  for STAs and  $O(s^3 \cdot |E|)$  for WSTAs, where  $|E|$  is the number of edges of the input instance. However, due to the compact representation of multiple edges, the size of the input instance may be much smaller than  $|E|$ .

Let  $\mathfrak{A}$  be an STA and  $\mathbf{I}$  an instance. Consider a node  $u$  of the input instance  $\mathbf{I}$  which has outgoing edges  $e_1, \dots, e_n$ . Suppose  $e_p, \dots, e_q$ , for some  $1 \leq p < q < n$ , have the same endpoint  $v$ , and  $e_{p-1}$  (if  $p > 1$ ) and  $e_{q+1}$  have an endpoint distinct from  $v$ . Then in the adjacency list of  $u$ , edges  $e_p, \dots, e_q$  are represented as a single edge  $e$  of multiplicity  $r = (q - p + 1)$ . Let  $w$  be the first child of  $v$  and  $f$  the edge from  $v$  to  $w$ . Let  $x$  be the endpoint of  $e_{q+1}$ , i.e., the next sibling of  $v$ . Suppose that during the computation of the reachable states we have already computed  $R =$

$reach(f)$  and  $R' = reach(e_{q+1})$ . To compute  $reach(e)$ , we not only have to simulate one step of the computation of the automaton  $\mathcal{A}$ , but actually a sequence of  $r$  steps to compute  $R_1 = reach(e_q), R_2 = reach(e_{q-1}), \dots, R_r = reach(e_p) = reach(e)$ . The crucial observation is that since there are only  $2^s$  distinct subsets of the state space, there are  $i < j \leq 2^{s+1}$  such that  $R_i = R_j$ . Of course this implies that for all  $k \geq i$  we have  $R_k = R_\ell$  for

$$\ell = ((k - i) \bmod (j - i)) + i.$$

Thus once we have found  $i$  and  $j$ , we can stop the simulation and compute  $R_r$  by simple arithmetic.

Instead of running this simulation every time it is needed, we can pre-compute a table that contains, for all subsets  $R$  and  $R'$ , the values  $i$  and  $j$  and the sets  $R_1, \dots, R_j$ . Such a table can be computed in time  $2^{3s}$ . Of course the table then simplifies the computation of the mappings  $reach$  and  $succ$ , because essentially we are dealing with a deterministic automaton now.  $\square$

### Acknowledgments

We thank an anonymous reviewer of LICS 2003 for pointing out to us that an unpublished part of Frank Neven's PhD thesis [23], pp. 128–130, defines the notion of *nondeterministic query automata* which precisely equals ours of STAs on (uncompressed) trees. Theorem 18 and Corollary 19 are proven there as well.

The third author's visit to LFCS, University of Edinburgh, was sponsored by Erwin Schrödinger Grant J2169 of the Austrian Research Fund (FWF).

### References

- [1] S. Abiteboul, P. Buneman, and D. Suciu. *Data on the Web*. Morgan Kaufmann Publishers, 2000.
- [2] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- [3] A. Brüggemann-Klein, M. Murata, and D. Wood. "Regular Tree and Regular Hedge Languages over Non-ranked Alphabets: Version 1, April 3, 2001". Technical Report HKUST-TCSC-2001-05, Hong Kong University of Science and Technology, Hong Kong SAR, China, 2001.
- [4] A. Brüggemann-Klein and D. Wood. "Caterpillars: A Context Specification Technique". *Markup Languages*, 2(1):81–106, 2000.
- [5] R. E. Bryant. "Graph-Based Algorithms for Boolean Function Manipulation". *IEEE Transactions on Computers*, C-35(8):677–691, Aug. 1986.
- [6] P. Buneman, M. Grohe, and C. Koch. "Path Queries on Compressed XML". In *Proceedings of the 29th International Conference on Very Large Data Bases (VLDB'03)*, pages 141–152, 2003.
- [7] J. Burch, E. Clarke, K. McMillan, D. Dill, and L. Hwang. "Symbolic Model Checking:  $10^{20}$  States and Beyond". In *Proceedings of the Annual IEEE Symposium on Logic in Computer Science (LICS)*, 1990.
- [8] E. M. Clarke, O. Grumberg, and D. Peled. *Model Checking*. MIT Press, 2000.
- [9] J. Doner. "Tree Acceptors and some of their Applications". *Journal of Computer and System Sciences*, 4:406–451, 1970.
- [10] R. Downey and M. Fellows. *Parameterized Complexity*. Springer-Verlag, 1999.
- [11] H.-D. Ebbinghaus and J. Flum. *Finite Model Theory*. Springer-Verlag, 1999. Second edition.
- [12] J. Flum, M. Frick, and M. Grohe. "Query Evaluation via Tree-Decompositions". In J. Van den Bussche and V. Vianu, editors, *Proc. of the 8th International Conference on Database Theory (ICDT'01)*, volume 1973 of *Lecture Notes in Computer Science*, pages 22–38, London, UK, Jan. 2001. Springer.
- [13] M. Frick and M. Grohe. "The Complexity of First-order and Monadic Second-order Logic Revisited". In *Proceedings of the 17th Annual IEEE Symposium on Logic in Computer Science (LICS)*, pages 215–224, 2002.
- [14] M. Frick, M. Grohe, and C. Koch. "Query Evaluation on Compressed Trees". In *Proceedings of the 18th Annual IEEE Symposium on Logic in Computer Science (LICS)*, Ottawa, Canada, June 2003.
- [15] G. Gottlob and C. Koch. "Monadic Queries over Tree-Structured Data". In *Proceedings of the 17th Annual IEEE Symposium on Logic in Computer Science (LICS)*, pages 189–202, Copenhagen, Denmark, July 2002.
- [16] G. Gottlob and C. Koch. "Monadic Datalog and the Expressive Power of Web Information Extraction Languages". *Journal of the ACM*, 51(1):1–40, Jan. 2004.
- [17] G. Gottlob, C. Koch, and R. Pichler. "Efficient Algorithms for Processing XPath Queries". In *Proceedings of the 28th International Conference on Very Large Data Bases (VLDB'02)*, Hong Kong, China, 2002.
- [18] G. Gottlob, C. Koch, and R. Pichler. "The Complexity of XPath Query Processing". In *Proceedings of the 22nd ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS'03)*, 2003.
- [19] T. J. Green, G. Miklau, M. Onizuka, and D. Suciu. "Processing XML Streams with Deterministic Automata". In *Proc. of the 9th International Conference on Database Theory (ICDT'03)*, 2003.
- [20] M. Grohe. "Parameterized Complexity for the Database Theorist". *SIGMOD Record*, 31(4), 2002.
- [21] H. Hosoya and B. C. Pierce. "Regular Expression Pattern Matching for XML". In *Proceedings of 28th Symposium on Principles of Programming Languages (POPL'01)*, pages 67–80. ACM Press, 2001.
- [22] T. Milo, D. Suciu, and V. Vianu. "Typechecking for XML Transformers". In *Proceedings of the ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS'00)*, pages 11–22, 2000.
- [23] F. Neven. *Design and Analysis of Query Languages for Structured Documents – A Formal and Logical Approach*. PhD thesis, Limburgs Universitair Centrum, 1999.

- [24] F. Neven. “Automata Theory for XML Researchers”. *SIGMOD Record*, **31**(3), Sept. 2002.
- [25] F. Neven and T. Schwentick. “Expressive and Efficient Pattern Languages for Tree-Structured Data”. In *Proceedings of the ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS’00)*, pages 145–156, Dallas, Texas, USA, 2000. ACM Press.
- [26] F. Neven and T. Schwentick. “Query Automata on Finite Trees”. *Theoretical Computer Science*, **275**:633–674, 2002.
- [27] L. Stockmeyer and A. Meyer. “Word Problems Requiring Exponential Time”. In *Proceedings of the 5th ACM Symposium on Theory of Computing*, pages 1–9, 1973.
- [28] J. Thatcher and J. Wright. “Generalized Finite Automata Theory with an Application to a Decision Problem of Second-order Logic”. *Mathematical Systems Theory*, **2**(1):57–81, 1968.
- [29] M. Y. Vardi. “The Complexity of Relational Query Languages”. In *Proc. 14th Annual ACM Symposium on Theory of Computing (STOC’82)*, pages 137–146, San Francisco, CA USA, May 1982.
- [30] World Wide Web Consortium. XML Path Language (XPath) Recommendation. <http://www.w3c.org/TR/xpath/>, Nov. 1999.