



**HAL**  
open science

# Querying Unranked Trees with Stepwise Tree Automata

Julien Carme, Joachim Niehren, Marc Tommasi

► **To cite this version:**

Julien Carme, Joachim Niehren, Marc Tommasi. Querying Unranked Trees with Stepwise Tree Automata. 19th International Conference on Rewriting Techniques and Applications, 2004, Aachen, Georgia. pp.105–118. inria-00536529

**HAL Id: inria-00536529**

**<https://inria.hal.science/inria-00536529>**

Submitted on 16 Nov 2010

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Querying Unranked Trees with Stepwise Tree Automata

Julien Carme      Joachim Niehren      Marc Tommasi

Mostrare project, INRIA Futurs, Lille, France

**Abstract.** The problem of selecting nodes in unranked trees is the most basic querying problem for XML. We propose *stepwise tree automata* for querying unranked trees. Stepwise tree automata can express the same monadic queries as monadic Datalog and monadic second-order logic. We prove this result by reduction to the ranked case, via a new systematic correspondence that relates unranked and ranked queries.

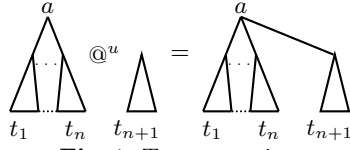
## 1 Introduction

Querying semi-structured documents is a base operation for information extraction from the Web or semi-structured databases. It requires expressive query languages whose queries can be answered efficiently [8]. The most widely known querying language these days is the W3C standard `XPath` (see e.g. [10, 9]).

Semi-structured documents in XML or HTML form *unranked trees* whose nodes may have an unbounded list of children. The most basic querying problem is to select sets of nodes in unranked trees. *Monadic queries* approach this problem declaratively. They specify sets of nodes in a tree that can then be computed by a generic algorithm.

We are interested in query languages that can describe all regular sets of nodes in trees. This property is satisfied by three classes of queries, those represented by tree automata [16, 12, 3, 13, 6], *monadic second-order logic* (MSO) [16, 8] and *monadic Datalog* [1, 7] over trees. Automata and Datalog queries can be answered in linear time. They are satisfactory in efficiency and expressiveness, in theory and practice [11].

Unranked trees are problematic in that they may be recursive in depth and breadth, in contrast to ranked trees. This additional level of recursion needs to be accounted for by recursive queries. In MSO and monadic Datalog, breadth recursion can be programmed from the `next_sibling` relation. Unfortunately, this relation cannot be expressed in  $WS\omega S$ , so that traditional results on ranked trees don't carry over for free. Selection automata [6] reduce breadth recursion to depth recursion, by operating on binary encodings of unranked trees. Encodings are problematic in that they alter locality and path properties; furthermore the close relationship to unranked Datalog queries gets lost. Hedge automata [16, 12, 3] express horizontal recursion by an extra recursion level in transition rules. This syntactic extension leads to numerous technical problems [13, 7] that one might prefer to avoid.



**Fig. 1.** Tree extension

In this paper, we propose *stepwise tree automata* for querying unranked trees. Stepwise tree automata are traditional tree automata that can either operate on unranked or ranked trees. They combine the advantages of selection and hedge automata. They model horizontal recursion by traversing siblings stepwise from the left to the right.

The algebraic approach behind stepwise tree automata yields a new systematic correspondence between queries for unranked and ranked trees. We elaborate this correspondence for monadic queries. We show that stepwise tree automata, monadic Datalog programs, and MSO can express the same monadic queries over unranked trees. We reduce this result to the case of ranked trees due to our new systematic correspondence. Specific proofs for unranked queries are not needed in contrast to [13, 7].

## 2 The algebras of unranked and ranked trees

An *unranked signature*  $\Sigma$  is a set of *symbol* ranged over by  $a, b$ . An ordered *unranked tree* or *u-tree*  $t$  over  $\Sigma$  satisfies the following abstract syntax:

$$t ::= a(t_1, \dots, t_n) \quad \text{where } n \geq 0.$$

Unordered unranked trees were investigated in [14, 15, 18]. We identify an unranked tree  $a()$  with the symbol  $a$ . We write  $\text{tree}^u$  for the set of unranked trees. The extension operator  $@^u : \text{tree}^u \times \text{tree}^u \rightarrow \text{tree}^u$  for unranked trees is depicted in Fig. 1. The extended tree  $t@^u t'$  is obtained from  $t$  by adjoining  $t'$  as next sibling of the last child of  $t$ :

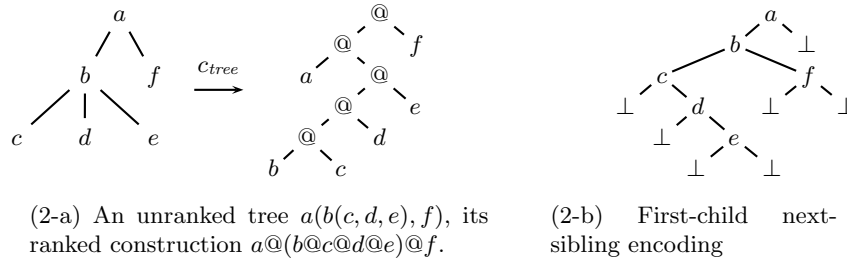
$$a(t_1, \dots, t_n)@^u t' = a(t_1, \dots, t_n, t')$$

Note that  $a(t_1, \dots, t_n) = a@^u t_1@^u \dots @^u t_n$  with parenthesis set from left to right. Tree extension  $@^u$  is neither associative nor commutative.

Let  $\Sigma_{@} = \Sigma \cup \{@\}$  be the *ranked signature* of function symbols with constants in  $\Sigma$  and a single binary function symbol  $@$ . Ranked trees over  $\Sigma_{@}$  are ground terms over  $\Sigma_{@}$ , i.e., binary trees that satisfy the grammar:

$$t ::= a \mid t_1 @ t_2$$

We omit parenthesis as in the  $\lambda$ -calculus; the ranked tree  $a@b@(c@b@a)$  for instance is  $(a@b)@((c@b)@a)$ . We write  $\text{tree}^r$  for the set of ranked trees over  $\Sigma_{@}$ .



Ranked trees can be constructed by the binary operator  $@^r : \mathbf{tree}^r \times \mathbf{tree}^r \rightarrow \mathbf{tree}^r$  which satisfies for all ranked trees  $t_1, t_2$ :

$$t_1 @^r t_2 = t_1 @ t_2$$

Unranked trees correspond precisely to ranked constructions with respect to the function  $c_{tree} : \mathbf{tree}^u \rightarrow \mathbf{tree}^r$  which satisfies for all unranked trees  $t_1, t_2$  and symbols  $a \in \Sigma$ :

$$c_{tree}(t_1 @^u t_2) = c_{tree}(t_1) @^r c_{tree}(t_2) \quad \text{and} \quad c_{tree}(a) = a$$

The idea of this binary construction is known from Curryng. An unranked tree describes an application of a function to a list of arguments. Its binary construction represents the Curried version of this function, receiving arguments one by one. We therefore write tree extension as function application  $@$ .

The set  $\mathbf{tree}^u$  of unranked trees over  $\Sigma$  with the extension operation  $@^u$  is a  $\Sigma_{@}$  algebra, as well as the set  $\mathbf{tree}^r$  of ranked trees over  $\Sigma_{@}$  with the operation  $@^r$ .

**Proposition 1.** *The construction function  $c_{tree} : \mathbf{tree}^u \rightarrow \mathbf{tree}^r$  is an isomorphism between  $\Sigma_{@}$ -algebras.*

Ranked and unranked trees thus have the same algebraic properties and the same finite automata [17, 5, 14].

Ranked constructions are binary representations of unranked trees. Previous approaches towards querying unranked trees rely on a different binary representation [8, 6], which encodes first-child and next-sibling relations. An example is given in Fig. 2-b. The new binary construction, however, permits to carry over traditional results from ranked to unranked trees more systematically.

### 3 Stepwise Tree Automata

*Stepwise tree automata*  $A$  over signature  $\Sigma$  are traditional tree automata ([4]) over the signature  $\Sigma_{@}$ . They consist of a finite set  $\mathbf{states}(A)$  of *states*, a set

$$\begin{aligned} \text{eval}_A^\alpha(a) &= \{q \mid a \rightarrow q \in \text{rules}(A)\} \\ \text{eval}_A^\alpha(t_1 @^\alpha t_2) &= \{q \mid q_1 \in \text{eval}_A^\alpha(t_1), q_2 \in \text{eval}_A^\alpha(t_2), q_1 @ q_2 \rightarrow q \in \text{rules}(A)\} \end{aligned}$$

**Fig. 3.** Ranked and unranked evaluation.

$\text{final}(A) \subseteq \text{states}(A)$  of *final states*, and a finite set of  $\text{rules}(A)$  of *transition rules* of two forms, where  $a \in \Sigma$  and  $q, q_1, q_2 \in \text{states}(A)$ :

$$a \rightarrow q \quad \text{or} \quad q_1 @ q_2 \rightarrow q$$

Stepwise tree automata can evaluate unranked and ranked trees, i.e.  $\alpha$ -trees where  $\alpha \in \{u, r\}$ . The  $\alpha$ -evaluator of  $A$  is the function  $\text{eval}_A^\alpha : \text{tree}^\alpha \rightarrow 2^{\text{states}(A)}$  defined in Fig. 3. Both evaluators only differ in the interpretation of the symbol  $@$ .

**Lemma 1.**  $\text{eval}_A^u(t) = \text{eval}_A^r(c_{\text{tree}}(t))$  for all unranked trees  $t$ .

Stepwise tree automata  $A$  recognize all  $\alpha$ -trees  $t$  that can be evaluated into a final state, i.e.,  $\text{eval}_A^\alpha(t) \cap \text{final}(A) \neq \emptyset$ . The  $\alpha$ -language  $L^\alpha(A)$  consists of all  $\alpha$ -trees recognized by  $A$ .

**Proposition 2.** *A stepwise tree automaton accepts an unranked tree if and only if it accepts its ranked construction, i.e., for all  $A$ :*

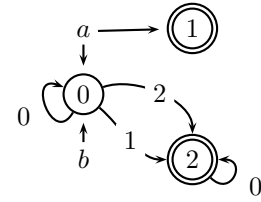
$$L^u(A) = c_{\text{tree}}^{-1}(L^r(A))$$

*Proof.* By Lemma 1 all unranked tree  $t$  satisfy:  $t \in L^u(A)$  iff  $\text{final}_A \cap \text{eval}_A^u(t) \neq \emptyset$  iff  $\text{final}_A \cap \text{eval}_A^r(c_{\text{tree}}(t)) \neq \emptyset$  iff  $c_{\text{tree}}(t) \in L^r(A)$ .

As a consequence, stepwise tree automata inherit numerous properties from traditional tree automata, even if interpreted over unranked trees. Recognizable unranked tree languages are closed under intersection, union, and complementation. Emptiness can be checked in linear time, and membership  $t \in L^u(A)$  in linear time  $O(|t| * |A|)$ .

**Example.** We define a stepwise tree automaton  $A$  with signature  $\Sigma = \{a, b\}$  that recognizes all unranked trees with at least one  $a$ -labeled leaf.

Automaton  $A$  is illustrated to the right. It has three states 0, 1, 2 two of which are final: 1, 2. A successful run of  $A$  on an unranked tree assigns state 1 in a non deterministic way to a unique  $a$ -leaf and labels by 2 all edges visited afterwards. All other nodes and edges are labeled by 0.



- An  $a$ -node can be the selected  $a$ -leaf:  $a \rightarrow 1$ .
- Any node can be assigned to state 0:  $a \rightarrow 0, b \rightarrow 0, 0@0 \rightarrow 0$ .
- The edge pointing to the selected  $a$ -leaf may go into state 2:  $0@1 \rightarrow 2$ .
- All edges visited later on may go into state 2, too:  $2@0 \rightarrow 2, 0@2 \rightarrow 2$ .

In Fig. 4, we show the unique successful run of  $A$  on the unranked tree  $a(b, a(a, b))$  and on its construction  $a@b@(a@a@b)$ . Both runs bisimulate each other. They evaluate the respective tree into the final state 2.

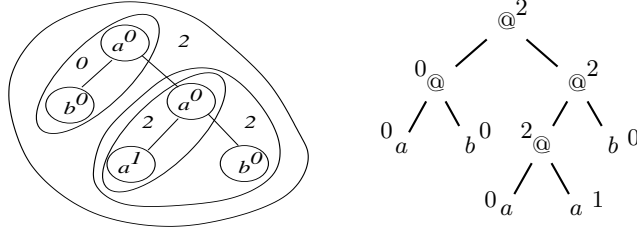


Fig. 4. An example run on an unranked tree and its construction

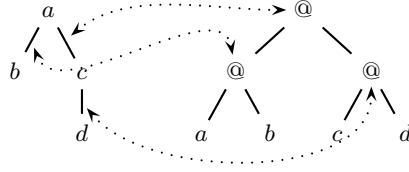


Fig. 5. Correspondence between edges and application nodes

#### 4 Monadic Queries

Queries on unranked trees should correspond precisely to their counterparts on ranked constructions. This requires a precise correspondence between the domains of unranked trees and their ranked constructions. The domain of a ranked tree  $t$  is the set of its nodes:

$$\begin{aligned} \text{dom}^r(t_1 @^r t_2) &= \text{dom}^r(t_1) \uplus \text{dom}^r(t_2) \uplus \{\text{root}\} \\ \text{dom}^r(a) &= \{\text{root}\} \end{aligned}$$

Disjoint union  $A \uplus B$  can be implemented by  $\{1\} \times A \cup \{2\} \times B$ . The ranked domain of  $a @ (a @ a)$  is then implemented by:

$$\{\text{root}, (1, \text{root}), (2, \text{root}), (2, (1, \text{root})), (2, (2, \text{root}))\}$$

As usual in mathematics, we will abstract from this implementation and talk about disjoint unions as if they were simple unions.

The nodes of unranked trees correspond to leaves in ranked constructions. But what do application nodes in construction trees correspond to? The example in Fig. 5 illustrates that they correspond precisely to edges of unranked trees. Every edge was added by some application step, and vice versa, every application step adds some edge. We therefore define the domain of an unranked tree as the union of its nodes and edges.

$$\begin{aligned} \text{dom}^u(t_1 @^u t_2) &= \text{dom}^u(t_1) \uplus \text{dom}^u(t_2) \uplus \{\text{last\_edge}\} \\ \text{dom}^u(a) &= \{\text{root}\} \end{aligned}$$

The `last_edge` of  $t_1 @^u t_2$  links the root of  $t_1$  to the root of  $t_2$ . Note that all nodes of  $t_1 @^u t_2$  either belong to  $t_1$  or  $t_2$ ; the root of  $t_1 @^u t_2$  is that of  $t_1$ .

$$c_{dom}(a)(\text{root}) = \text{root}, \quad c_{dom}(t_1 @^u t_2)(\pi) = \begin{cases} \text{root} & \text{if } \pi = \text{last\_edge} \\ c_{dom}(t_1)(\pi) & \text{if } \pi \in \text{dom}^u(t_1) \\ c_{dom}(t_2)(\pi) & \text{if } \pi \in \text{dom}^u(t_2) \end{cases}$$

**Fig. 6.** Definition of the correspondence on domains  $c_{dom}$ .

The *correspondence*  $c_{dom}(t)$  for an unranked tree  $t$  is a function between the domains of  $t$  and its ranked construction defined in Fig. 6 by recursion over the construction of  $t$ :

$$c_{dom}(t) : \text{dom}^u(t) \rightarrow \text{dom}^r(c_{tree}(t))$$

**Definition 1.** A monadic query is a function  $q$  that maps trees to subsets of their domain. This definition applies to ranked and unranked trees, i.e., for both  $\alpha \in \{u, r\}$ . A monadic  $\alpha$ -query  $q$  satisfies for all  $t \in \text{tree}^\alpha$ :

$$q(t) \subseteq \text{dom}^\alpha(t)$$

We restrict ourselves to monadic queries; more general n-ary queries map to n-tuples of elements of the domain [2]. Monadic queries over unranked and ranked trees correspond. An unranked monadic query  $q$  corresponds to the ranked monadic queries  $c_{query}(q)$  which satisfies for all  $t \in \text{tree}^r$ :

$$c_{query}(q)(t) = c_{dom}(t)(q(c_{tree}^{-1}(t)))$$

All previous query notions for unranked trees [8, 6] only talk about nodes. Our extension with edges, however, is necessary to keep the symmetry to ranked queries, the reason for the simplicity of our approach.

## 5 Automata Queries

In the remainder of the paper, we will discuss regular monadic queries. These can be defined by tree automata, monadic second-order logic, and monadic Datalog. Here, we start with tree automata.

We next consider monadic queries as tree languages over the alphabet  $\Sigma \times \text{Bool}$ . Such languages were already studied in [17]. Given a ranked tree  $t$  over the signature  $\Sigma \times \text{Bool}$  and  $i \in \{1, 2\}$ , let  $\text{proj}^i(t)$  be the ranked tree obtained by projecting all labels in  $t$  to their  $i$ 'th component. For monadic queries  $q$  and ranked trees  $t$  let  $\text{zip}(t, q)$  be the ranked tree over the extended signature  $\Sigma \times \text{Bool}$  with  $\text{proj}^1(\text{zip}(t, q)) = t$  and  $\text{proj}^2(\text{zip}(t, q)) = q$ .

**Definition 2.** A monadic query  $q$  for ranked trees is regular if the set  $\{\text{zip}(t, q) \mid t \in \text{tree}^r\}$  can be recognized by a tree automaton. A monadic query  $q$  for unranked trees is regular if  $c_{query}(q)$  is.

This traditional definition of regular monadic queries is simple for ranked trees but has drawbacks otherwise. First, it is not obvious how to compute such

$$\frac{a \rightarrow r(\text{root}) \in \text{rules}(A) \quad r \in \text{runs}_A^\alpha(a)}{r \in \text{runs}_A^\alpha(a)} \quad \frac{r|_{\text{dom}^\alpha(t_1)} \in \text{runs}_A^\alpha(t_1) \quad r|_{\text{dom}^\alpha(t_2)} \in \text{runs}_A^\alpha(t_2) \quad r(\text{head}^\alpha(t_1))@r(\text{head}^\alpha(t_2)) \rightarrow r(\text{head}(t_1@^\alpha t_2)) \in \text{rules}(A)}{r \in \text{runs}_A^\alpha(t_1@^\alpha t_2)}$$

**Fig. 7.** Runs  $\text{runs}_A^\alpha(t)$  of stepwise tree automata  $A$  on  $\alpha$ -trees  $t$ .

queries efficiently, second, it is not obvious how to express them in monadic Datalog, and third, the definition of regular unranked queries depends on the correspondence to ranked queries.

Run-based queries [15, 6] with stepwise tree automata resolve these problems. We define them parametrically for ranked and unranked trees. Runs of stepwise tree automaton on  $\alpha$ -trees  $t$  associate states to all elements of the domain of  $t$ . Sequences of children in unranked trees are visited *stepwise* from the left to the right, while annotating edges to the children by states. See Fig. 4 for an example. More formally, let the *head* of a ranked tree be its *root*; the head of a non-constant unranked tree is *last\_edge*, and the head of a constant unranked tree the *root*:

$$\text{head}^r(t) = \text{root}, \quad \text{head}^u(t_1@^u t_2) = \text{last\_edge}, \quad \text{head}^u(a) = \text{root}.$$

A *run* of a tree automaton  $A$  on an  $\alpha$ -tree  $t$  is a function labeling elements of the domain of  $t$  by states of  $A$ :

$$r : \text{dom}^\alpha(t) \rightarrow \text{states}(A)$$

such that all transitions are licensed by rules of  $A$ . If  $t = t_1@^u t_2$  then the restrictions of  $r$  to the domains of  $t_1$  and  $t_2$  must be runs and the annotation of the head of  $t$  must be justified. Furthermore, annotations of constants must be licensed. These conditions are captured by the inference rules in Fig. 7.

**Lemma 2.** *Let  $A$  be a stepwise tree automaton and  $t$  an  $\alpha$ -tree, then:*

$$\text{eval}_A^\alpha(t) = \{r(\text{head}^\alpha(t)) \mid r \in \text{runs}_A^\alpha(t)\}$$

*Proof.* By induction on the construction of unranked trees. If  $t = a$  then  $\text{eval}^\alpha(a) = \{q \mid a \rightarrow q \in \text{rules}(A)\} = \{r(\text{root}) \mid r \in \text{runs}_A^\alpha(a)\}$ . For  $t = t_1@^\alpha t_2$  we have  $\text{eval}^\alpha(t) = \{q \mid q_i \in \text{eval}_A^\alpha(t_i), q_1@q_2 \rightarrow q \in \text{rules}(A)\}$  which is equal to  $\{q \mid r_i \in \text{runs}_A^\alpha(t_i), r_1(\text{head}^\alpha(t_1))@r_2(\text{head}^\alpha(t_2)) \rightarrow q \in \text{rules}(A)\}$  by induction hypothesis. The definition of runs in Fig. 7 yields  $\{r(\text{head}_i^\alpha(t)) \mid r \in \text{runs}_A^\alpha(t)\}$ .

A run  $r$  of an automaton  $A$  on an  $\alpha$ -tree  $t$  is *successful* if  $r(\text{head}^\alpha(t)) \in \text{final}(A)$ . Let  $\text{succ\_runs}_A^\alpha(t)$  be the set of successful runs of  $A$  on  $t \in \text{tree}^\alpha$ .

**Definition 3.** *A pair of a tree automaton  $A$  and a set  $Q \subseteq \text{states}(A)$  defines a monadic query which selects all elements from the domain of a tree  $t$  that are labeled by a state in  $Q$  in some successful run of  $A$  on  $t$ :*

$$\text{query}_{A,Q}^\alpha(t) = \{\pi \in \text{dom}^\alpha(t) \mid r \in \text{succ\_runs}_A^\alpha(t), r(\pi) \in Q\}$$



Selection automata [6] similarly express queries for unranked trees, but rely on universal quantification over successful runs, and use a binary encoding of unranked trees in contrast to the definition above.

**Example.** Reconsider the automaton  $A$  from Section 3:  $\text{query}_{A,\{1\}}^u$  defines the set of all  $a$ -leaves in unranked trees. Note that no automaton query with a bottom-up deterministic automaton can compute the same query, since it couldn't distinguish different  $a$ -nodes.

**Proposition 3.** *A monadic query on ranked trees is regular if and only if it is equal to some  $\text{query}_{A,Q}^r$ .*

The proof relies on two standard automata transformations. The idea of the transformation from regular to run based queries  $\text{query}_{A,Q}^r$  is memorize the Boolean values in  $\text{zip}(t, q)$  in automata states. In order to generalize Proposition 3 to unranked trees, we establish the correspondence between unranked and ranked run-based queries.

**Theorem 1.** *Queries with stepwise tree automata on ranked and unranked trees correspond:*

$$\text{query}_{A,Q}^r = c_{\text{query}}(\text{query}_{A,Q}^u)$$

*Proof.* We first note that runs of stepwise automata on unranked trees and ranked constructions correspond. For all  $t \in \text{tree}^r$ , we can prove:

$$\text{runs}_A^u(c_{\text{tree}}^{-1}(t)) = \{r \circ c_{\text{dom}}(t) \mid r \in \text{runs}_A^r(t)\}$$

The theorem follows from straightforward calculations. For all  $t \in \text{tree}^r$ :

$$\begin{aligned} c_{\text{query}}(\text{query}_{A,Q}^u)(t) &= c_{\text{dom}}(t)(\text{query}_{A,Q}^u(c_{\text{tree}}^{-1}(t))) \\ &= \{\pi \mid r \in \text{runs}_A^u(c_{\text{tree}}^{-1}(t)), r(c_{\text{dom}}(t)^{-1}(\pi)) \in Q\} \\ &= \{\pi \mid r' \in \text{runs}_A^r(t), r'(c_{\text{dom}}(t)(c_{\text{dom}}(t)^{-1}(\pi))) \in Q\} \\ &= \text{query}_{A,Q}^r(t) \quad \square \end{aligned}$$

## 6 Monadic Second Order Logic

We next represent regular monadic queries in ranked and unranked trees in monadic second-order logic (MSO).

The domain of the logical structure induced by an  $\alpha$ -tree  $t$  is  $\text{dom}^\alpha(t)$ . The signature  $R^r$  for structures of ranked trees contains the following relation symbols:

$$R^r = \{\text{child}_1, \text{child}_2, \text{root}, \text{leaf}\} \cup \{\text{label}_a \mid a \in \Sigma_\circ\}$$

The binary relations  $\text{child}_1$  and  $\text{child}_2$  relate nodes to their first resp. second child. Unary relations  $\text{label}_a$  hold for all nodes labeled by  $a \in \Sigma$ . Furthermore, we permit the unary relations  $\text{root}$  and  $\text{leaf}$ .

Logical structures for unranked trees have the following signature:

$$R^u = \{\text{first\_edge}, \text{next\_edge}, \text{target}, \text{root}, \text{leaf}\} \cup \{\text{label}_a \mid a \in \Sigma_\circ\}$$

The binary relation `first_edge` holds between a node and the edge to its first child. The `next_edge` relation links an edge with target  $\pi$  to the edge whose target is the next sibling of  $\pi$ . The `target` relation holds between an edge and its target node.

Let  $x, y, z$  range over an infinite set `Vars` of node variables and  $p, q$  over an infinite set `Preds` of monadic predicates. The logics  $\text{MSO}^\alpha$  have the following formulas:

$$\phi ::= B_n(x_1, \dots, x_n) \mid p(x) \mid \phi \wedge \phi' \mid \neg\phi \mid \exists x\phi \mid \exists p\phi$$

where  $B_n \in R^\alpha$  is a predicate with fixed tree interpretation of arity  $n$ . Note that the relations `root` and `leaf` could be expressed by the remaining relations in  $\text{MSO}^\alpha$ . We add them anyway, as they will be needed in monadic Datalog later on.

Let  $\phi$  be an  $\text{MSO}^\alpha$ -formula,  $t$  an  $\alpha$ -tree, and  $\sigma$  an assignment of variables into the domain of  $t$  and of predicates into the powerset of this domain. We write  $t, \sigma \models_{\text{MSO}^\alpha} \phi$  if  $\phi$  becomes true in  $t$  under  $\sigma$ . Every formula  $\phi(x)$  with a single variable  $x$  defines monadic query:

$$\text{query}_{\phi(x)}^\alpha(t) = \{\sigma(x) \mid t, \sigma \models_{\text{MSO}^\alpha} \phi\}$$

**Theorem 2.** [Thatcher & Wright [17]] *Monadic queries expressed in monadic second order logic over ranked trees  $\text{MSO}^r$  are regular.*

**Theorem 3.** *Ranked and unranked monadic queries expressed in monadic second-order logic correspond, and are thus regular; corresponding queries can be computed in linear time:*

$$\{\text{query}_{\phi(x)}^r \mid \phi \in \text{MSO}^r\} = \{c_{\text{query}}(\text{query}_{\phi'(x)}^u) \mid \phi' \in \text{MSO}^u\}$$

Fig. 8 presents forth and back translations between  $\text{MSO}^u$  and  $\text{MSO}^r$ . We have to show for every  $\phi \in \text{MSO}^u$  that  $c_{\text{query}}(\text{query}_{\phi(x)}^u) = \text{query}_{\llbracket \phi(x) \rrbracket_r}^r$ , and the analogous property for the back translation. We proceed by structural induction over formulas. The base cases contains the difficulty, the induction step being straightforward.

We sketch the proof for formula `first_edge`( $x, y$ ) to illustrate the principles. Consider an  $u$ -tree  $t$  and a variable assignment  $\sigma$  under which `first_edge`( $x, y$ ) becomes true. There exists a  $u$ -tree  $t_0 = t_1 @ t_2$  involved in the construction of  $t$  such that  $\sigma(x)$  is the root of  $t_1$  and  $\sigma(y)$  is the edge from  $t_1$  to  $t_2$ . Since this is the first edge,  $t_1$  is a constant. Therefore, in the corresponding ranked tree, the node  $c_{\text{dom}}(\sigma(x))$  is a leaf and we have  $\text{child}_1(c_{\text{dom}}(\sigma(y)), c_{\text{dom}}(\sigma(x)))$ . The converse is proved in a similar way.

The case of `target`( $x, y$ ) is more tedious as it relies on the recursive `lar`( $x, y$ ) formula, stating that  $x$  is a leaf, whose last ancestor to the right is  $y$ . This means that  $y$  denotes the up most node with  $\text{child}_1^*(y, x)$ . A model of `target`( $x, y$ ) and its translation  $\llbracket \phi(x) \rrbracket_r$  in Fig.10.

Auxiliary predicates:

$$\begin{aligned} \text{ar}'(x, p) &=_{\text{def}} p(x) \wedge \forall y \forall z ((\text{child}_1(y, z) \wedge p(z)) \rightarrow p(y)) \\ \text{ar}(x, p) &=_{\text{def}} \text{ar}'(x, p) \wedge \forall p' (\text{ar}'(x, p') \rightarrow \text{subset}(p, p')) \\ \text{lar}(x, y) &=_{\text{def}} \text{leaf}(x) \wedge \exists p. p(y) \wedge \text{ar}(x, p) \wedge (\text{root}(y) \vee \exists y' \text{child}_2(y', y)) \end{aligned}$$

Logical connectives

$$\begin{aligned} \llbracket \exists x \psi \rrbracket_r &=_{\text{def}} \exists x \llbracket \psi \rrbracket_r \\ \llbracket \exists p \psi \rrbracket_r &=_{\text{def}} \exists p \llbracket \psi \rrbracket_r \\ \llbracket \psi \wedge \psi' \rrbracket_r &=_{\text{def}} \llbracket \psi \rrbracket_r \wedge \llbracket \psi' \rrbracket_r \\ \llbracket \neg \psi \rrbracket_r &=_{\text{def}} \neg \llbracket \psi \rrbracket_r \\ \llbracket p(x) \rrbracket_r &=_{\text{def}} p(x) \end{aligned}$$

Node relations

$$\begin{aligned} \llbracket \text{first\_edge}(x, y) \rrbracket_r &=_{\text{def}} \text{child}_1(y, x) \wedge \text{leaf}(x) \\ \llbracket \text{next\_edge}(x, y) \rrbracket_r &=_{\text{def}} \exists z (\text{child}_1(y, x) \wedge \text{child}_1(x, z)) \\ \llbracket \text{target}(x, y) \rrbracket_r &=_{\text{def}} \exists z (\text{child}_2(x, z) \wedge \text{lar}(y, z)) \\ \llbracket \text{label}_a(x) \rrbracket_r &=_{\text{def}} \text{label}_a(x) \quad (a \in \Sigma) \\ \llbracket \text{last\_edge}(x, y) \rrbracket_r &=_{\text{def}} \exists z (\text{lar}(x, y) \wedge \text{child}_1(y, z)) \\ \llbracket \text{root}(x) \rrbracket_r &=_{\text{def}} \exists y (\text{lar}(x, y) \wedge \text{root}(y)) \\ \llbracket \text{leaf}(x) \rrbracket_r &=_{\text{def}} \exists y \text{child}_2(y, x) \wedge \text{leaf}(x) \end{aligned}$$

Fig. 8. Unranked into ranked MSO

$$\begin{aligned} \llbracket \exists x \psi \rrbracket_u &=_{\text{def}} \exists x \llbracket \psi \rrbracket_u & \llbracket \text{child}_1(x, y) \rrbracket_u &=_{\text{def}} \text{first\_edge}(y, x) \vee \text{next\_edge}(y, x) \\ \llbracket \exists p \psi \rrbracket_u &=_{\text{def}} \exists p \llbracket \psi \rrbracket_u & \llbracket \text{child}_2(x, y) \rrbracket_u &=_{\text{def}} (\text{target}(x, y) \wedge \text{leaf}(y)) \\ \llbracket \psi \wedge \psi' \rrbracket_u &=_{\text{def}} \llbracket \psi \rrbracket_u \wedge \llbracket \psi' \rrbracket_u & & \vee \exists z (\text{target}(x, z) \wedge \text{last\_edge}(z, y)) \\ \llbracket \neg \psi \rrbracket_u &=_{\text{def}} \neg \llbracket \psi \rrbracket_u & \llbracket \text{label}_a(x) \rrbracket_u &=_{\text{def}} \text{label}_a(x) \quad (a \in \Sigma) \\ \llbracket p(x) \rrbracket_u &=_{\text{def}} p(x) & \llbracket \text{root}(x) \rrbracket_u &=_{\text{def}} (\text{leaf}(x) \wedge \text{root}(x)) \\ & & & \vee \exists z (\text{root}(z) \wedge \text{last\_edge}(z, x)) \\ & & \llbracket \text{leaf}(x) \rrbracket_u &=_{\text{def}} \text{root}(x) \vee \exists z \text{target}(z, x) \end{aligned}$$

Fig. 9. Ranked into unranked MSO

## 7 Monadic Datalog

We next express regular monadic queries in Monadic Datalog, a logic programming language, and discuss the expressive power compared to automata and MSO queries, for ranked and unranked trees.

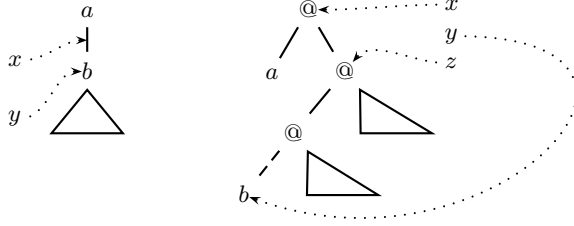
We consider monadic Datalog in trees without negation. The languages  $\text{Datalog}^\alpha$  have the same signatures as  $\text{MSO}^\alpha$ . The programs of  $\text{Datalog}^\alpha$  are logic program without function symbols, predefined n-ary predicates in  $R^\alpha$ , and free monadic predicates  $p, q \in \text{Preds}$ . More precisely, a program  $P \in \text{Datalog}^\alpha$  is a finite set of rules of the form:

$$p(x) :- \text{Body}$$

where  $\text{Body}$  is a sequence of goals with n-ary predicates  $B_n \in R^\alpha$ :

$$\text{Body} ::= B_n(x_1, \dots, x_n) \mid p(x) \mid \text{Body}, \text{Body}'$$

Every program of  $\text{Datalog}^\alpha$  can be seen as a formula of  $\text{MSO}^\alpha$ ; sets of clauses are conjunctions, clauses  $p(x) :- \text{Body}$  are universally quantified implications  $\forall \text{Vars}. p(x) \leftarrow \text{Body}$ , and bodies  $\text{Body}$  conjunctions of goals.



**Fig. 10.** A solution of  $\text{target}(x, y)$  on the left; a solution of its ranked translation  $\llbracket \text{target}(x, y) \rrbracket_r = \exists z (\text{child}_2(x, z) \wedge \text{lar}(z, y))$  on the right.

We interpret programs in  $\text{Datalog}^\alpha$  in the least fixed point semantics over the structures of  $\text{MSO}^\alpha$ . For every program  $P \in \text{Datalog}^\alpha$ , predicate  $p \in \text{Preds}$ , and  $t \in \text{tree}^\alpha$  let  $T_{P,t}^\omega(p)$  be the least solution of  $P$  over the tree structure of the  $\alpha$ -tree  $t$  for predicate  $p$ . This yields a notion of monadic queries:

$$\text{query}_{P(p)}^\alpha(t) = T_{P,t}^\omega(p)$$

Least fixed points can be expressed in  $\text{MSO}$ . As a consequence, every query in  $\text{Datalog}^\alpha$  can be expressed in linear time in  $\text{MSO}^\alpha$ . This shows that ranked queries in  $\text{Datalog}^r$  are regular (Theorem 2).

Ranked monadic run-based automata queries can always be expressed in ranked monadic  $\text{Datalog}$ ; the reduction is in linear time. The resulting  $\text{Datalog}$  program models the two phases of the linear time algorithm for answering automata queries: the first bottom up phase computes all states of all nodes seen in all runs of the automaton, and a top down phase selects all nodes labeled by selection states in successful runs.

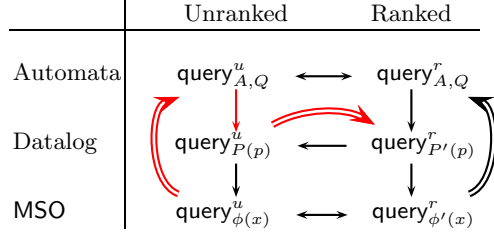
**Theorem 4.** *Ranked and unranked monadic queries expressed in monadic  $\text{Datalog}$  correspond, and are thus regular:*

$$\{\text{query}_{P(p)}^r \mid P \in \text{Datalog}^r\} = \{c_{\text{query}}(\text{query}_{P'(p)}^u) \mid P' \in \text{Datalog}^u\}$$

*Corresponding unranked queries can be computed from ranked queries in linear time; the converse is not true.*

It suffices to encode ranked  $\text{Datalog}$  queries into corresponding unranked queries in linear time. The translation basically refines the encoding from  $\text{MSO}^r$  into  $\text{MSO}^u$ : roughly, a rule  $p(x) :- \text{Body}$  is translated into  $p(x) :- \llbracket \text{Body} \rrbracket_u$ . Conjunctions in  $\llbracket \text{Body} \rrbracket_u$  can be replaced by commas, existential quantifications can be omitted, i.e., replaced by implicit universal quantification in rules. Disjunctions as in the definitions of  $\llbracket \text{child}_1(x, y) \rrbracket_u$ ,  $\llbracket \text{child}_2(x, y) \rrbracket_u$ , and  $\llbracket \text{root}(x) \rrbracket_u$  can be expressed by multiple rules. Such a rewriting, however, spoils linear time. We can circumvent this problem following Gottlob and Koch [7]: we normalise programs of  $\text{Datalog}^r$  into *tree marking normal form* (TMNF) in linear time before translation. TMNF programs have of forms:

$$\begin{array}{ll} p(x) :- B_1(x). & p(x) :- q(y), B_2(y, x). \\ p(x) :- p_0(x), p_1(x). & p(x) :- q(y), B_2(x, y). \end{array}$$



**Fig. 11.** Summary of reductions. Solid lines are in linear time, double lines are non elementary. Black lines are proved, red lines induced.

$$\begin{array}{c}
 t = a(t_1, \dots, t_n) \qquad \forall 1 \leq i \leq n : r|_{\text{nodes}(t_i)} \in \text{run}_H(t_i) \\
 a(L) \rightarrow r(\text{root}(t)) \in \text{rules}(H) \qquad r(\text{root}(t_1)) \dots r(\text{root}(t_n)) \in L \\
 \hline
 r \in \text{run}_H(t)
 \end{array}$$

**Fig. 12.** Runs of Hedge Automata

where  $B_n$  is a  $n$ -ary predicate of  $R^r$ . On ranked TMNF programs the reduction  $\llbracket - \rrbracket_u$  can clearly be done in linear time.

The inverse translation can be composed from our translations so far, which are summarized in Fig. 11. We first reduce unranked Datalog queries to  $\text{MSO}^u$ , then to  $\text{MSO}^r$ , move to ranked automata queries and then to ranked Datalog queries. The overall reduction has nonelementary complexity.

Note that we cannot specialize the translation  $\llbracket - \rrbracket_r$  from  $\text{MSO}^u$  to  $\text{MSO}^r$  into a translation from  $\text{Datalog}^u$  to  $\text{Datalog}^r$ . The problem is that we cannot express the auxiliary binary predicate  $\text{lar}(x, y)$ . Its definition is recursive, but only monadic recursive predicates can be defined in monadic Datalog.

**Corollary 1.** *Stepwise tree automata, monadic Datalog, and MSO capture the class of regular monadic queries over unranked trees.*

This is a corollary of our correspondences between ranked and unranked queries and traditional result on ranked trees. All unranked automata queries  $\text{query}_{A,Q}^u$  can be expressed in linear time in  $\text{Datalog}^u$ , by indirection over ranked automata queries and  $\text{Datalog}^r$ . Unranked queries in  $\text{MSO}^u$  can be expressed by unranked automata queries by reduction to the ranked case.

## 8 Hedge Automata

We finally show how to express monadic queries with hedge automata [16, 3] in linear time with stepwise tree automata. A hedge automaton  $H$  over  $\Sigma$  consists of a set  $\text{states}(H)$  of states, a set  $\text{final}(H)$  of final states, and set set of transition rules of the form  $a(L) \rightarrow q$  where  $L$  is a regular set of words over states.

Runs of hedge automata  $H$  on unranked trees  $t$  are functions  $r : \text{nodes}(t) \rightarrow \text{states}(H)$  that satisfy the inference rule in Fig 12. A hedge automaton  $H$  and a

$$\begin{aligned}
\text{states}(\text{step}(H)) &= \text{states}(H) \uplus \biguplus_{a \in \Sigma, q \in \text{states}(H)} \text{states}(H_{a,q}) \\
\text{rules}(\text{step}(H)) &= \bigcup_{a \in \Sigma, q \in \text{states}(H)} \text{rules}(H_{a,q}) \\
&\quad \cup \{p \xrightarrow{\epsilon} q \mid p \in \text{final}(H_{a,q}), q \in \text{states}(H), a \in \Sigma\} \\
&\quad \cup \{a \rightarrow p \mid p \in \text{init}(H_{a,q}), q \in \text{states}(H)\} \\
\text{final}(\text{step}(H)) &= \text{final}(H)
\end{aligned}$$

**Fig. 13.** Hedge automata into stepwise tree automata

set of states  $Q \subseteq \text{states}(H)$  defines a monadic query for unranked trees:

$$\text{query}_{H,Q}(t) = \{\pi \in \text{nodes}(t) \mid r \in \text{succ\_runs}_H(t), r(\pi) \in Q\}$$

For translating hedge automata into stepwise tree automata, we need to represent all regular language  $L$  in transition rules explicitly. We use a sequence of finite word automata  $(H_{a,q})_{a \in \Sigma, q \in \text{states}(H)}$  over the alphabet  $\text{states}(H)$  to do so.

**Proposition 4.** *Queries by hedge automata can be translated in linear time to queries by stepwise tree automata.*

*Proof.* Given an hedge automaton  $H$  we define a stepwise tree automaton  $\text{step}(H)$  by unifying all subautomata  $H_{a,q}$  into a single finite automaton. We then add all states of  $\text{states}(A)$  to this automaton, and link them to all final states  $p$  of  $H_{a,q}$  through  $\epsilon$ -transitions  $p \xrightarrow{\epsilon} q$ . We add rules  $a \rightarrow q'$  for all initial states  $q'$  of some  $H_{a,q}$ . The final states of  $\text{step}(H)$  are those in  $\text{final}(H)$ , not those in  $\text{final}(H_{a,q})$ . The complete construction is detailed in Fig. 13. It remains to show that every run of  $H$  can be simulated by a run of  $\text{step}(H)$ .  $\square$

**Acknowledgments:** We would like to thank our colleagues from the Mostrare project for their continuous support, particularly R. Gilleron and S. Tison, and the anonymous referees for hinting us towards Currification.

## References

1. R. Baumgartner, S. Flesca, and G. Gottlob. Visual web information extraction with Lixto. In *The Very Large Data Bases Journal*, pages 119–128, 2001.
2. A. Berlea and H. Seidl. Binary queries. In *Proceedings of Extreme Markup Languages*, Montreal, 2002.
3. A. Brüggemann-Klein, M. Murata, and D. Wood. Regular tree and regular hedge languages over unranked alphabets. Technical report, 2001.
4. H. Comon, M. Dauchet, R. Gilleron, F. Jacquemard, D. Lugiez, S. Tison, and M. Tommasi. Tree automata techniques and applications. Online book, 450 pages. Available at: <http://www.grappa.univ-lille3.fr/tata>, 1997.
5. B. Courcelle. On recognizable sets and tree automata. In H. Ait-Kaci and M. Nivat, editors, *Resolution of Equations in Algebraic Structures, Algebraic Techniques*, volume 1, chapter 3, pages 93–126. Academic Press, 1989.
6. M. Frick, M. Grohe, and C. Koch. Query evaluation on compressed trees. In *Proceedings of the IEEE Symposium on Logic In Computer Sciences*, Ottawa, 2003.

7. G. Gottlob and C. Koch. Monadic datalog and the expressive power of languages for web information extraction. In *Proceedings of the ACM Symposium on Principle of Databases Systems*, pages 17–28, 2002.
8. G. Gottlob and C. Koch. Monadic queries over tree-structured data. In *Proceedings of the IEEE Symposium on Logic In Computer Sciences*, Copenhagen, 2002.
9. G. Gottlob, C. Koch, and R. Pichler. The complexity of XPATH query evaluation. In *Proceedings of the ACM Symposium on Principle of Databases Systems*, pages 179–190. 2003.
10. G. Gottlob, C. Koch, and R. Pichler. XPATH processing in a nutshell. *ACM SIGMOD Record*, 32(2):21–27, 2003.
11. C. Koch. Efficient processing of expressive node-selecting queries on XML data in secondary storage: A tree automata-based approach. In *Proceedings of the International Conference on Very Large Data Bases*, 2003.
12. A. Neumann, H. Seidl. Locating matches of tree patterns in forests. In *Foundations of Software Technology and Theoretical Computer Science*, pages 134–145, 1998.
13. F. Neven and T. Schwentick. Query automata over finite trees. *Theoretical Computer Science*, 275(1-2):633–674, 2002.
14. J. Niehren and A. Podelski. Feature automata and recognizable sets of feature trees. In *Proceedings of TAPSOFT'93*, volume 668 of *LNCS*, pages 356–375, 1993.
15. H. Seidl, T. Schwentick, and A. Muscholl. Numerical document queries. In *Proc. of the IEEE Symposium on Principles of Database Systems*, pages 155–166, 2003.
16. J. W. Thatcher. Characterizing derivation trees of context-free grammars through a generalization of automata theory. *J. of Comp. and Syst. Sci.*, 1:317–322, 1967.
17. J. W. Thatcher, J. B. Wright. Generalized finite automata with an application to a decision problem of second-order logic. *Math. System Theory*, 2:57–82, 1968.
18. S. D. Zilio and D. Lugiez. XML schema, tree logic and sheaves automata. In R. Nieuwenhuis, editor, *Proceedings of RTA*, volume 2706 of *LNCS*, pages 246–263. 2003.