



**HAL**  
open science

# Efficient Inclusion Checking for Deterministic Tree Automata and XML Schemas

Jérôme Champavère, Rémi Gilleron, Aurélien Lemay, Joachim Niehren

► **To cite this version:**

Jérôme Champavère, Rémi Gilleron, Aurélien Lemay, Joachim Niehren. Efficient Inclusion Checking for Deterministic Tree Automata and XML Schemas. *Information and Computation*, 2009, 207 (11), pp.1181-1208. 10.1016/j.ic.2009.03.003 . inria-00366082v3

**HAL Id: inria-00366082**

**<https://inria.hal.science/inria-00366082v3>**

Submitted on 10 Oct 2009

**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.

# Efficient Inclusion Checking for Deterministic Tree Automata and XML Schemas

Jérôme Champavère<sup>\*1,2</sup>, Rémi Gilleron<sup>1,2</sup>, Aurélien Lemay<sup>1,2</sup>,  
Joachim Niehren<sup>2</sup>

*Inria Lille – Nord Europe, Parc scientifique de la Haute Borne, 40 avenue Halley,  
Bât. A, Park Plaza, 59650 Villeneuve d'Ascq, France*

---

## Abstract

We present algorithms for testing language inclusion  $L(A) \subseteq L(B)$  between tree automata in time  $O(|A| \cdot |B|)$  where  $B$  is deterministic (bottom-up or top-down). We extend our algorithms for testing inclusion of automata for unranked trees  $A$  in deterministic DTDs or deterministic EDTDs with restrained competition  $D$  in time  $O(|A| \cdot |\Sigma| \cdot |D|)$ . Previous algorithms were less efficient or less general.

*Key words:* tree automata, language inclusion, algorithmic complexity, XML schemas

---

## 1. Introduction

Language inclusion for tree automata is a basic decision problem that is closely related to universality and equivalence [1, 2, 3]. Tree automata algorithms are generally relevant for XML document processing [4, 5, 6, 7]. Regarding inclusion, typical applications are inverse type checking for tree transducers [8] and schema-guided query induction [9]. The latter was the motivation for the present study. There, candidate queries produced by the learning process are to be checked for consistency with deterministic DTDs, such as for HTML.

We investigate language inclusion  $L(A) \subseteq L(B)$  for tree automata  $A$  and  $B$  under the assumption that  $B$  is bottom-up deterministic or top-down deterministic, not necessarily  $A$ . Without this assumption, the problem be-

---

\*Corresponding author.

*Email address:* [jerome.champavere@lifl.fr](mailto:jerome.champavere@lifl.fr) (Jérôme Champavère)

<sup>1</sup>Laboratoire d'Informatique Fondamentale de Lille, Université de Lille.

<sup>2</sup>Mostrare project (<http://mostrare.futurs.inria.fr/>), Inria Lille–Nord Europe.

comes DEXPTIME-complete [3]. Deterministic language inclusion still subsumes universality of deterministic tree automata  $L(B) = T_\Sigma$  up to a linear time reduction, as well as equivalence of languages of two deterministic automata  $L(A) = L(B)$ . Conversely, one can reduce inclusion to equivalence in PTIME, since  $L(A) \subseteq L(B)$  if and only if  $L(A) \cap L(B) = L(A)$ . However, this leads to a reduction in quadratic time  $O(|A| \cdot |B|)$ , so we cannot rely on equivalence testing (as by comparing cardinalities [2] or unique minimal deterministic automata) for efficient inclusion testing.

The well-known naive test for inclusion in bottom-up deterministic tree automata for ranked trees goes through complementation. It first computes an automaton  $B^c$  that recognizes the complement of the language of  $B$ , and then checks whether the intersection automaton for  $B^c$  and  $A$  has a non-empty language. The problematic step is the completion of  $B$  before complementing its final states, since completion might require adding rules for all possible left-hand sides. The overall running time may thus become  $O(|A| \cdot |\Sigma| \cdot |B|^n)$ , which is exponential in the maximal rank  $n$  of function symbols in the signature  $\Sigma$ . This time complexity can be reduced by turning the maximal arity of function symbols in ranked trees to 2. It is folklore that one can transform ranked trees into binary trees, and automata correspondingly. The problem here is to preserve bottom-up determinism, while the size of automata must remain linear. We can solve this problem by using curried encodings of unranked tree into binary trees, as proposed for stepwise tree automata [10, 1]. Thereby we obtain an inclusion test for the ranked case in time  $O(|A| \cdot |\Sigma| \cdot |B|^2)$ . This is still too much in practice with XML schemas, where  $A$  and  $B$  may be of size 500 and  $\Sigma$  of size 100; these orders of magnitude can be observed, e.g., in the DTDs of the corpus studied by Bex *et al.* [11].

Our first contribution is a more efficient algorithm that test inclusion in bottom-up deterministic tree automata in time  $O(|A| \cdot |B|)$ . This bound is independent of the size of the signature  $\Sigma$ , even if it is not fixed. We establish our algorithm for stepwise tree automata over binary trees in the first step, and then lift it *via* currying to standard tree automata for ranked trees over arbitrary signatures and to stepwise tree automata over unranked trees.

As a second contribution, we show how to test language inclusion of stepwise tree automata  $A$  for unranked trees in deterministic DTDs  $D$  in time  $O(|A| \cdot |\Sigma| \cdot |D|)$ . Determinism for DTDs is required by the XML standards. Our algorithm first computes Glushkov automata for all regular expressions in  $D$  in time  $O(|\Sigma| \cdot |D|)$ . This is possible since we assume  $D$  to be deterministic DTDs [12]. The second step is more tedious. We would like to transform the whole collection of Glushkov automata into a single bottom-up deterministic stepwise tree automaton of the same size. Unfortunately, this seems

difficult to achieve, since the usual construction of Martens & Niehren [13] eliminates  $\epsilon$ -rules on the fly, which may lead to a quadratic blowup of the number of rules (not the number of states).

We solve this problem by introducing bottom-up deterministic *factorized tree automata*. These are tree automata with  $\epsilon$ -rules, which represent deterministic stepwise tree automata more compactly, and in particular the collection of Glushkov automata of a DTD of linear size. Factorized tree automata have two sorts of states, which play the roles of hedge and tree states in alternative automata notions for unranked trees [14, 4]. The difficulty is to define the appropriate notion of determinism for factorized tree automata, and to adapt the inclusion test to the case where  $B$  is a deterministic factorized tree automaton.

Our results can be applied if  $A$  is a hedge automaton [1, 15, 16], with finite word automata for horizontal languages, since such hedge automata can be translated in linear time to stepwise tree automata. Note, however, that the notion of (bottom-up) determinism for hedge automata is unsatisfactory [13] so that we cannot choose  $B$  to be a deterministic hedge automaton even if the horizontal language is defined by a deterministic finite word automaton.

The situation becomes slightly different if  $A$  is a tree automaton recognizing firstchild-nextsibling encodings of unranked trees, and  $D$  a DTD. The problem is that the conversion of  $A$  into a stepwise tree automaton may lead to a quadratic size increase. In this case, however, we can encode DTDs into top-down deterministic tree automata that recognize firstchild-nextsibling encodings of unranked trees, and reduce the inclusion problem to the case of inclusion in deterministic finite word automata. This yields a worst case running time of  $O(|A| \cdot |\Sigma| \cdot |D|)$ , too. As we show, the same algorithm applies if  $D$  is a deterministic extended DTD (EDTD) with restrained competition [17]. These were introduced in order to reason about schema definitions in the W3C standard XML Schema [18, 19], and relaxations thereof.

Our algorithm generalizes the inclusion test of Martens, Neven, Schwentick & Bex [18] (see Section 10 of the reference), where  $A$  and  $D$  are both limited to deterministic EDTDs with restrained competition. The presentation of our algorithm differs in that we rely on inclusion in top-down deterministic tree automata *via* firstchild-nextsibling encoding as an intermediate step, while they reduce the problem to inclusion in deterministic finite word automata directly (*via* the main theorem of this article). Furthermore, we provide a precise complexity analysis for the first time (which is not fully obvious).

*Related Work.* Our new algorithm for testing inclusion in bottom-up deterministic factorized tree automata  $B$  is relevant for schemas defined in Relax NG [20]. This holds for those definitions that can be made bottom-up

deterministic without combinatorial explosion. Furthermore, we can permit arbitrary Relax NG schemas as automata  $A$  on the left, where no determinism is required.

The folklore algorithms for testing inclusion in top-down deterministic automata (by reduction to finite automata for path languages) lead us to a generalization of the inclusion test for deterministic restrained competition EDTDs [18]. This is useful for testing inclusion of Relax NG in XML Schema for instance.

Compared to the conference version of the present article at LATA'08 [21], we have added new results on early failure detection, incrementality, experiments, and complete proofs. Furthermore, we added the alternative algorithm for inclusion in top-down deterministic automata and in restrained competition EDTDs. We simplified the presentation of our algorithms in many places. Meanwhile, the inclusion test presented here has been integrated into a system for schema-guided query induction [9], where it proves its efficiency in practice.

The complexity of inclusion for various fragments of DTDs and EDTDs was first studied by Martens, Neven & Schwentick [22]. They assume the same types of language definitions on both sides. When applied to deterministic DTDs, the same complexity results seem obtainable when refining the efficiency analysis provided there. In any case, our algorithm permits richer left-hand sides (as needed in schema-guided query induction) without increasing in complexity.

Heuristic algorithms for inclusion between non-deterministic automata and applications that avoid the high worst-case complexity were proposed by Tozawa & Hagiya [23]. Even though motivated by XML Schema, for which better algorithms are available meanwhile (due to top-down determinism), they are relevant for Relax NG where no notion of determinism is imposed *a priori*.

*Outline.* In Section 2, we reduce inclusion for ranked tree automata to the binary case. An efficient incremental algorithm for binary tree automata is given in Section 3. In Section 4, we introduce deterministic factorized tree automata and lift the inclusion test. In Section 5 we apply it to test inclusion of automata in deterministic DTDs. Section 6 presents experimental results. Section 7 studies inclusion in top-down deterministic tree automata and restrained competition extended DTDs. Appendix A details the implementation.

## 2. Standard Tree Automata for Ranked Trees

We reduce the inclusion problem of tree automata for ranked trees [1] to the case of binary trees with a single binary function symbol.

A ranked signature  $\Sigma$  is a finite set of function symbols  $f \in \Sigma$ , each of which has an arity  $ar(f) \geq 0$ . A constant  $a \in \Sigma$  is a function symbol of arity 0. A tree  $t \in T_\Sigma$  is either a constant  $a \in \Sigma$  or a tuple  $f(t_1, \dots, t_n)$  consisting of a function symbol  $f$  with  $ar(f) = n$  and trees  $t_1, \dots, t_n \in T_\Sigma$ .

A *tree automaton*  $A$  over  $\Sigma$  with  $\epsilon$ -rules consists of a finite set  $sta(A)$  of states, a subset  $fin(A) \subseteq sta(A)$  of final states, and a set  $rul(A) \subseteq sta(A)^2 \uplus (\cup_{n \geq 0} \{f \in \Sigma \mid ar(f) = n\} \times sta(A)^{n+1})$ . We denote such rules as  $p' \xrightarrow{\epsilon} p$  or  $f(p_1, \dots, p_n) \rightarrow p$ , where  $f \in \Sigma$  has arity  $n$  and  $p_1, \dots, p_n, p, p' \in sta(A)$ . Furthermore, we write  $p' \xrightarrow{\epsilon}_A p$  iff  $p' \xrightarrow{\epsilon} p \in rul(A)$ ,  $\xrightarrow{\epsilon}_A^*$  for the reflexive transitive closure of  $\xrightarrow{\epsilon}_A$ , and  $\xrightarrow{\epsilon}_A^{\leq 1}$  for the union of  $\xrightarrow{\epsilon}_A$  and the identity relation on  $sta(A)$ .

The *size* of an  $\epsilon$ -rule  $p \xrightarrow{\epsilon} p'$  is 2 and that of a rule  $f(p_1, \dots, p_n) \rightarrow p$  is  $n+2$ . The *size*  $|A|$  of  $A$  is the cardinality of  $sta(A)$ , denoted  $|sta(A)|$ , plus the sum of the sizes of the rules of  $A$ , which we denote  $|rul(A)|$ . The cardinality  $|\Sigma|$  of the signature  $\Sigma$  is ignored, since it is irrelevant for algorithms that care only about used symbols. Every tree automaton  $A$  defines an evaluator  $eval_A : T_{\Sigma \cup sta(A)} \rightarrow 2^{sta(A)}$  such that:

$$eval_A(f(t_1, \dots, t_n)) = \left\{ p \mid \begin{array}{l} p_1 \in eval_A(t_1), \dots, p_n \in eval_A(t_n), \\ f(p_1, \dots, p_n) \rightarrow p' \in rul(A), p' \xrightarrow{\epsilon}_A p \end{array} \right\}$$

and  $eval_A(p) = \{p\}$ . A tree  $t \in T_\Sigma$  is *accepted* by  $A$  if  $fin(A) \cap eval_A(t) \neq \emptyset$ . The *language*  $L(A)$  is the set of trees accepted by  $A$ .

A tree automaton is (bottom-up) *deterministic* if it has no  $\epsilon$ -rules, and if no two rules have the same left-hand side. It is complete if there are rules for all potential left-hand sides. It is well-known that deterministic complete tree automata can be complemented in linear time, by switching the final states.

*Deterministic inclusion.* We will study the deterministic inclusion problem for tree automata. Its input consists of a ranked signature  $\Sigma$ , a possibly non-deterministic tree automaton  $A$  with  $\epsilon$ -rules, and a deterministic tree automaton  $B$ , both with signature  $\Sigma$ , and its output is the truth value of  $L(A) \subseteq L(B)$ .

We can deal with this problem by restriction to stepwise signatures  $\Sigma_{@}$ , which consist of a single binary function symbol  $@$  and a finite set of constants  $a \in \Sigma$ . A stepwise tree automaton over binary trees is a tree automaton over a stepwise signature [10]. We use the infix notation in automata rules and thus write  $q_1 @ q_2 \rightarrow q$  instead of  $@(q_1, q_2) \rightarrow q$ .

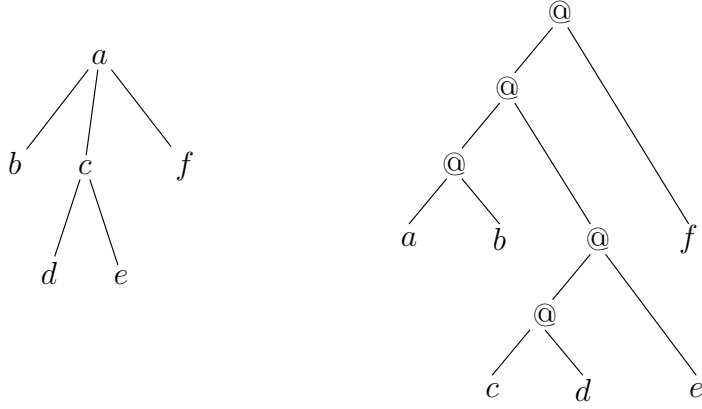


Figure 1: Currying the ranked tree  $a(b, c(d, e), f)$  into the binary tree  $a@b@(c@d@e)@f$ .

In Section 5, we will see how to interpret stepwise tree automata over unranked trees *via* binary encoding. Here, we use the same binary encoding for interpretation over ranked trees with arbitrary signatures.

**Proposition 1.** *The deterministic inclusion problem for standard tree automata over ranked trees can be reduced in linear time to the deterministic inclusion problem for stepwise tree automata over binary trees.*

We first encode ranked trees into binary trees *via* currying. Given a ranked signature  $\Sigma$  we define the corresponding signature  $\Sigma_{@} = \{@\} \uplus \Sigma$  whereby all symbols of  $\Sigma$  become constants. We use infix notation for the binary symbol  $@$  and write  $t_1@t_2$  instead of  $@(t_1, t_2)$ . Furthermore, we assume that omitted parentheses have priority to the left, i.e., we write  $t_1@t_2@t_3$  instead of  $(t_1@t_2)@t_3$ . Currying is defined by a function  $curry : T_{\Sigma} \rightarrow T_{\Sigma_{@}}$  which for all trees  $t_1, \dots, t_n \in T_{\Sigma}$  and  $f \in \Sigma$  satisfies:

$$curry(f(t_1, \dots, t_n)) = f@curry(t_1)@\dots@curry(t_n)$$

For instance,  $a(b, c(d, e), f)$  is mapped to  $a@b@(c@d@e)@f$ , which is the infix notation for the tree  $@(@(@(a, b), @(c, d), e), f)$ , as shown in Figure 1.

Now we encode tree automata  $A$  over  $\Sigma$  into stepwise tree automata  $step(A)$  over  $\Sigma_{@}$ , such that the language is preserved up to currying, i.e., such that  $L(step(A)) = curry(L(A))$ . The states of  $step(A)$  are the prefixes of left-hand sides of rules in  $A$ , i.e., words in  $\Sigma(sta(A))^*$ :

$$sta(step(A)) = \{fq_1 \dots q_i \mid f(q_1, \dots, q_n) \rightarrow q \in rul(A), 0 \leq i \leq n\} \uplus sta(A)$$

$$\begin{array}{c}
\frac{f(q_1, \dots, q_n) \rightarrow q \in \text{rul}(A) \quad 1 \leq i < n}{f q_1 \dots q_{i-1} @_{q_i} \rightarrow f q_1 \dots q_i \in \text{rul}(\text{step}(A))} \\
\frac{f q_1 \dots q_{i-1} @_{q_i} \rightarrow f q_1 \dots q_i \in \text{rul}(\text{step}(A))}{f q_1 \dots q_{n-1} @_{q_n} \rightarrow q \in \text{rul}(\text{step}(A))}
\end{array}
\qquad
\frac{a \rightarrow q \in \text{rul}(A)}{a \rightarrow q \in \text{rul}(\text{step}(A))}$$

Figure 2: Transforming ranked tree automata into stepwise tree automata.

The rules of  $\text{step}(A)$  are given in Figure 2. They extend prefixes step by step by states  $q_i$  according to the rules of  $A$ . Since constants cannot be extended, we need to distinguish two cases.

**Lemma 2.** *The encoding of tree automata  $A$  over  $\Sigma$  into stepwise tree automata  $\text{step}(A)$  over  $\Sigma_{@}$  preserves determinism, the tree language modulo currying, and the automata size up to a constant factor of 3.*

As a consequence,  $L(A) \subseteq L(B)$  is equivalent to  $L(\text{step}(A)) \subseteq L(\text{step}(B))$ , and can be tested in this way modulo a linear time transformation. Most importantly, the determinism of  $B$  carries over to  $\text{step}(B)$ .

### 3. Stepwise Tree Automata for Binary Trees

We present our new algorithm for testing deterministic inclusion in the case of stepwise tree automata over binary trees. We start with a characterization of deterministic inclusion, express it by a Datalog program [24, 25], and then turn it into an efficient algorithm, which is non-trivial.

#### 3.1. Ground Datalog

For the sake of self-containedness, we recall folklore results on ground Datalog (see, e.g., Gottlob *et al.* [26]). A ground Datalog program is set of Horn clauses, without function symbols, variables, and negation. More formally, it is build from a ranked signature  $\Gamma$  with constants  $c \in \Gamma$  and predicates  $\mathbf{p} \in \Gamma$ , each of which has an arity  $\text{ar}(\mathbf{p}) \geq 0$ . A literal is a term of the form  $\mathbf{p}(c_1, \dots, c_{\text{ar}(\mathbf{p})})$ . We write  $\text{lit}(\Gamma)$  for the set of all literals over  $\Gamma$ . A (Horn) clause, written as “ $L :- L_1, \dots, L_k$ .”, is a pair in  $\text{lit}(\Gamma) \times \text{lit}(\Gamma)^k$ , where  $k \geq 0$ . As usual, we write “ $L$ .” instead of “ $L :- .$ ”, where  $k = 0$ . A ground Datalog program  $P$  over  $\Gamma$  is a finite set of Horn clauses over  $\Gamma$ . The size  $|P|$  of a Datalog program  $P$  is the overall number of occurrences of symbols in its clauses.

Every ground Datalog program  $P$  over  $\Gamma$  has a unique least fixed point  $\text{lfp}(P)$  (since there is no negation). This is the least set of literals over  $\Gamma$  that satisfies for all  $L :- L_1, \dots, L_k$ . in  $P$ , that  $L_1 \dots, L_k \in \text{lfp}(P)$  implies



$L \in \text{lfp}(P)$ . Least fixed points are always finite sets (in the absence of function symbols).

The next theorem states that least fixed points can be computed efficiently (since there are no variables).

**Theorem 3 (Efficiency of Ground Datalog).** *For every signature  $\Gamma$  and ground Datalog program  $P$  over  $\Gamma$ , the least fixed point  $\text{lfp}(P)$  can be computed in linear time  $O(|P|)$ .*

Note that the upper bound  $O(|P|)$  may depend on the arities of predicates in  $P$ , but not on the arities of the other predicates of  $\Gamma$ .

*Proof.* A program  $P$  defines a hypergraph, whose edges are the tuples  $(L, L_1, \dots, L_k)$  with  $L :- L_1, \dots, L_k$  in  $P$ . The least fixed point  $\text{lfp}(P)$  is the set of literals accessible in this hypergraph. It is well-known that accessible components of graphs can be computed in linear time. The same holds for hypergraphs, under the often implicit condition that  $L = L'$  can be tested in time  $O(1)$ . This condition is clearly valid if all predicates in  $P$  have arity 0. For this case, the theorem is folklore (see, e.g., Minoux [27]).

For an arbitrary signature  $\Gamma$ , we consider  $\text{lit}(\Gamma)$  as the ranked signature without constants, in which all literals become predicates of arity 0. Literals  $L$  over  $\Gamma$  corresponds one-to-one to literals  $L()$  over  $\text{lit}(\Gamma)$ . Thus, every ground Datalog program over  $\Gamma$  can be transformed into a ground Datalog program over  $\text{lit}(\Gamma)$  in time  $O(|P|)$ . In order to test  $L = L'$  in time  $O(1)$  as required above, we have to replace all literals in  $P$  by numbers, such that different occurrences of the same literal are mapped to the same number. This can be done in time  $O(|P|)$  by using a prefix tree, that memorizes all the numbers assigned to literals seen so far. For instance, the prefix-tree  $\text{lit}(\mathbf{p}_1(c_1(c_2(1), c_3(2)), \mathbf{p}_0(c_1(3))))$  memorizes the assignments of  $\mathbf{p}_1(c_1, c_2)$  to 1, of  $\mathbf{p}_1(c_1, c_3)$  to 2, and of  $\mathbf{p}_0(c_1)$  to 3.  $\square$

We can refine least fixed points from sets to multisets, by counting for every literal the multiplicity with which it can be added to the least fixed point, where  $\#S$  denotes the cardinality of the set  $S$ :

$$\begin{aligned} \text{lfp}^\#(P) &: \text{lit}(\Gamma) \rightarrow \mathbb{N} \cup \{0\} \\ \text{lfp}^\#(P)(L) &= \#\{R \in \text{lfp}(P)^k \mid L :- R. \text{ in } P, k \geq 0\} \end{aligned}$$

Note that  $L \in \text{lfp}(P)$  if and only if  $\text{lfp}^\#(P)(L) > 0$  by definition. The next corollary shows that multiplicities of literals in least fixed points can be computed efficiently.

**Corollary 4.** *For every signature  $\Gamma$  and ground Datalog program  $P$  over  $\Gamma$ , a representation of the least fixed point with multiplicities  $\text{lfp}^\#(P)$  can be computed in linear time  $O(|P|)$ .*

---

$\frac{a \rightarrow p \in \text{rul}(A)}{\text{acc}(p)}$	$\frac{p_1 @ p_2 \rightarrow p \in \text{rul}(A)}{\text{acc}(p) :- \text{acc}(p_1), \text{acc}(p_2)}$	$\frac{p' \xrightarrow{\epsilon}_A p}{\text{acc}(p) :- \text{acc}(p')}$
$\frac{p \in \text{fin}(A)}{\text{coacc}(p)}$	$\frac{p_1 @ p_2 \rightarrow p \in \text{rul}(A)}{\text{coacc}(p_1) :- \text{coacc}(p), \text{acc}(p_2)}$	$\frac{p' \xrightarrow{\epsilon}_A p}{\text{coacc}(p') :- \text{coacc}(p)}$

---

Figure 3: Accessible and co-accessible states of  $A$ .

We represent  $\text{lfp}^\#(P)$  by its restriction to  $\text{lfp}(P)$ , i.e., by the relation  $\{(L, \text{lfp}^\#(P)(L)) \mid L \in \text{lfp}(P)\}$  which contains all non-zero values of  $\text{lfp}^\#(P)$ . Thereby, we avoid enumerating the elements of the complement of the least fixed point  $\text{lfp}(P)^c$ .

*Proof.* The relation  $\{(L, \text{lfp}^\#(P)(L)) \mid L \in \text{lfp}(P)\}$  can be computed from  $P$  and  $\text{lfp}(P)$  in time  $O(|P|)$ , by inspecting all clauses of  $P$  exactly once, and counting for all literals how often they appear on the left-hand side of clauses whose literals on the right-hand side all belong to  $\text{lfp}(P)$ . It is thus sufficient to compute the set  $\text{lfp}(P)$  in time  $O(|P|)$ . This can be done by Theorem 3.  $\square$

### 3.2. Characterization of Inclusion

Let  $A$  be a tree automaton over  $\Sigma$ . We call a state  $p \in \text{sta}(A)$  *accessible* (or sometimes *reachable*) if there exists a tree  $t \in T_\Sigma$  such that  $p \in \text{eval}_A(t)$ , and *co-accessible* if there exists a tree  $C[p] \in T_{\Sigma \cup \{p\}}$  with a unique occurrence of  $p$  (i.e., a context with hole marker  $p$ ) such that  $\text{eval}_A(C[p]) \cap \text{fin}(A) \neq \emptyset$ . For every term  $s \in T_\Sigma$ , we denote by  $C[s] \in T_\Sigma$  the term obtained by replacing the unique occurrence of  $p$  in  $C[p]$  by  $s$ .

We call  $A$  *productive* if all its states are accessible and co-accessible. Note that productive automata may become unproductive by completion, since sink states are not co-accessible. The ground Datalog program in Figure 3 computes all accessible and co-accessible states of an automaton  $A$ . The rules of  $A$  are transformed to clauses of the Datalog program. The overall number of such clauses is linear in the size of  $A$ , so the least fixed point can be computed in linear time by Theorem 3. States that are unaccessible or not co-accessible, and all rules using them can be safely removed from  $A$ . This renders  $A$  productive in linear time, while preserving its language.

We will use the following notion for stepwise tree automata  $A$  with states  $p_1, p_2 \in \text{sta}(A)$ , meaning that evaluation may proceed in this pair:

$$A \models p_1 @ p_2 \Leftrightarrow_{df} \exists p \in \text{sta}(A). p_1 @ p_2 \rightarrow p \in \text{rul}(A)$$

Let  $B$  be another automaton over  $\Sigma$  but without  $\epsilon$ -rules. The *product*  $A \times B$  has state set  $sta(A) \times sta(B)$ , and rules inferred as follows:

$$\frac{a \rightarrow p \in rul(A)}{a \rightarrow (p, q)} \quad \frac{p_1 @ p_2 \rightarrow p \in rul(A) \quad q_1 @ q_2 \rightarrow q \in rul(B)}{(p_1, q_1) @ (p_2, q_2) \rightarrow (p, q)} \quad \frac{p' \xrightarrow{A} p \quad q \in sta(B)}{(p', q) \xrightarrow{\epsilon} (p, q)}$$

We do not care about final states of  $A \times B$  since these are useless in our characterization of inclusion. The following property of states  $p$  of  $A$  and  $q$  of  $B$  is equivalent to accessibility of the pair  $(p, q)$  in  $A \times B$ :

$$A, B \models \text{acc}(p, q) \Leftrightarrow_{df} (p, q) \text{ accessible in } A \times B$$

Language inclusion  $L(A) \subseteq L(B)$  fails under the following three conditions:

$A, B \models \text{fail}_0$ : there exists a rule  $a \rightarrow p \in rul(A)$  but no state  $q \in sta(B)$  such that  $a \rightarrow q \in rul(B)$ ;

$A, B \models \text{fail}_1$ : there exist states  $p_1, p_2, q_1, q_2$  such that  $A, B \models \text{acc}(p_1, q_1)$ ,  $A, B \models \text{acc}(p_2, q_2)$ ,  $A \models p_1 @ p_2$  and  $B \not\models q_1 @ q_2$ ;

$A, B \models \text{fail}_2$ : there exist  $p \in fin(A)$  and  $q \notin fin(B)$  such that  $A, B \models \text{acc}(p, q)$ .

We compose the properties of automata pairs by first-order connectives: we write  $A, B \models \phi_1 \vee \phi_2$  iff  $A, B \models \phi_1$  or  $A, B \models \phi_2$ , and similarly for the other first-order connectives such as  $A, B \models \phi \Rightarrow \phi'$ .

**Proposition 5.** *Inclusion  $L(A) \subseteq L(B)$  for productive stepwise tree automata  $A$  with  $\epsilon$ -rules and deterministic stepwise tree automata  $B$  fails iff  $A, B \models \text{fail}_0 \vee \text{fail}_1 \vee \text{fail}_2$*

*Proof.* For soundness, we suppose that one of the failure conditions holds, and show that some tree  $t \in L(A)$  witnesses inclusion failure, i.e.,  $t \notin L(B)$ .

$A, B \models \text{fail}_0$ . Let us consider a rule  $a \rightarrow p \in rul(A)$  such that no rule  $a \rightarrow q \in rul(B)$  exists. Since  $A$  is productive, state  $p$  is co-accessible, i.e., there exists a term  $C[p] \in T_{\Sigma \cup \{p\}}$  with a single occurrence of  $p$  such that  $eval_A(C[p]) \cap fin(A) \neq \emptyset$ . Hence  $C[a] \in L(A)$ . But  $C[a] \notin L(B)$  because there is no rule  $a \rightarrow q \in rul(B)$ .

$A, B \models \text{fail}_1$ . There exists  $t_1 \in T_\Sigma$  such that  $(p_1, q_1) \in eval_{A \times B}(t_1)$  by accessibility of  $(p_1, q_1)$  and there exists  $t_2 \in T_\Sigma$  such that  $(p_2, q_2) \in eval_{A \times B}(t_2)$  by accessibility of  $(p_2, q_2)$ . Since

$p_1 @ p_2 \rightarrow p \in \text{rul}(A)$  we also get  $p \in \text{eval}_A(t_1 @ t_2)$  by definition of  $\text{eval}_A$ . Furthermore since  $A$  is productive there exists a term  $C[p] \in T_{\Sigma \cup \{p\}}$  with a single occurrence  $p$  such that  $C[t_1 @ t_2] \in L(A)$ . Since  $B$  is deterministic it follows that  $q_1 \in \text{eval}_B(t_1)$  and  $q_2 \in \text{eval}_B(t_2)$  are unique. By hypothesis there is no  $q$  such that  $q_1 @ q_2 \rightarrow q \in \text{rul}(B)$ , so that  $C[t_1 @ t_2] \notin L(B)$ .

$A, B \models \text{fail}_2$ . There are  $p \in \text{fin}(A)$  and  $q \notin \text{fin}(B)$  such that  $(p, q)$  is accessible. Thus, there exists  $t \in T_\Sigma$  such that  $(p, q) \in \text{eval}_{A \times B}(t)$ . The state  $p$  is final in  $A$ , hence  $t \in L(A)$ . Since  $B$  is deterministic  $q \in \text{eval}_B(t)$  is unique but  $q$  not final in  $B$  implies  $t \notin L(B)$ .

For completeness, we assume that there exists a tree  $t \in L(A)$  such that  $t \notin L(B)$ , and show that some failure condition holds. There are two cases to be considered, depending on  $\text{eval}_B(t)$ .

- (i) Assume  $\text{eval}_B(t) = \emptyset$ . There exists a minimal subtree  $t'$  of  $t$  such that  $\text{eval}_B(t') = \emptyset$ , too. If  $t' = a$  is a leaf then  $\text{eval}_A(a) \neq \emptyset$ , since  $t \in L(A)$ , and  $\text{eval}_B(a) = \emptyset$ , hence  $A, B \models \text{fail}_0$ . If  $t' = t_1 @ t_2$ , then there exist  $p_1 \in \text{eval}_A(t_1)$ ,  $p_2 \in \text{eval}_A(t_2)$  and  $p_1 @ p_2 \rightarrow p \in \text{rul}(A)$ , since  $t \in L(A)$ . Since  $t'$  is defined as a minimal subtree and  $B$  is deterministic,  $\text{eval}_B(t_1) = \{q_1\}$ ,  $\text{eval}_B(t_2) = \{q_2\}$ , and since  $\text{eval}_B(t') = \emptyset$ , there is no rule  $q_1 @ q_2 \rightarrow q \in \text{rul}(B)$ . This shows  $A, B \models \text{fail}_1$ .
- (ii) If  $\text{eval}_B(t) \neq \emptyset$  then there exists  $q \in \text{eval}_B(t)$ ; since  $B$  is deterministic,  $q$  is necessarily unique. Since  $t \notin L(B)$  this yields  $q \notin \text{fin}(B)$ . Moreover, since  $t \in L(A)$ , there exists  $p \in \text{eval}_A(t) \cap \text{fin}(A)$ . Thus,  $A, B \models \text{fail}_2$  holds.  $\square$

### 3.3. Testing Characterization in Ground Datalog

We next transform stepwise tree automata  $A$  and  $B$  into a ground Datalog program by which to test the failure conditions.

Figure 4 presents the transformation of two automata  $A, B$  into a ground Datalog program  $D_0(A, B)$ , which tests whether  $A, B \models \text{fail}_0$  or  $A, B \models \text{fail}_2$ . The clauses produced from  $A$  and  $B$  by three transformation rules ( $\text{acc}_{/1}$ ), ( $\text{acc}_{/2}$ ), and ( $\text{acc}_{/3}$ ) compute the accessibility relation  $\text{acc}$  of  $A \times B$  as usual. Clearly,  $\text{acc}(p, q) \in \text{lfp}(D_0(A, B))$  iff  $A, B \models \text{acc}(p, q)$ . The Datalog clauses produced by transformation rules ( $\text{fail}_0$ ) and ( $\text{fail}_2$ ) serve for computing the predicates  $\text{fail}_0$  and  $\text{fail}_2$ . By construction,  $\text{fail}_0 \in \text{lfp}(D_0(A, B))$  iff  $A, B \models \text{fail}_0$  and  $\text{fail}_2 \in \text{lfp}(D_0(A, B))$  iff  $A, B \models \text{fail}_2$ .

---


$$\begin{array}{l}
(\text{acc}/_1) \frac{a \rightarrow p \in \text{rul}(A) \quad a \rightarrow q \in \text{rul}(B)}{\text{acc}(p, q)}. \\
(\text{acc}/_2) \frac{p' \xrightarrow{\epsilon}_A p \in \text{rul}(A) \quad q \in \text{sta}(B)}{\text{acc}(p, q) :- \text{acc}(p', q)}. \\
(\text{acc}/_3) \frac{p_1 @ p_2 \rightarrow p \in \text{rul}(A) \quad q_1 @ q_2 \rightarrow q \in \text{rul}(B)}{\text{acc}(p, q) :- \text{acc}(p_1, q_1), \text{acc}(p_2, q_2)}. \\
(\text{fail}_0) \frac{a \rightarrow p \in \text{rul}(A) \quad \nexists q \in \text{sta}(B). a \rightarrow q \in \text{rul}(B)}{\text{fail}_0}. \\
(\text{fail}_2) \frac{p \in \text{fn}(A) \quad q \notin \text{fn}(B)}{\text{fail}_2 :- \text{acc}(p, q)}.
\end{array}$$


---

Figure 4: Datalog program  $D_0(A, B)$  testing  $A, B \models \text{fail}_0$  and  $A, B \models \text{fail}_2$ .

---


$$\begin{array}{l}
(\text{frb}/_1) \frac{A \models p_1 @ p_2 \quad B \not\models q_1 @ q_2}{\text{frb}(p_2, q_2) :- \text{acc}(p_1, q_1)}. \\
(\text{frb}/_2) \frac{A \models p_1 @ p_2 \quad B \not\models q_1 @ q_2}{\text{frb}(p_1, q_1) :- \text{acc}(p_2, q_2)}. \\
(\text{fail}_1) \frac{p \in \text{sta}(A) \quad q \in \text{sta}(B)}{\text{fail}_1 :- \text{acc}(p, q), \text{frb}(p, q)}.
\end{array}$$


---

Figure 5: Datalog program  $D_1(A, B)$  testing  $A, B \models \text{fail}_1$ .

Datalog program  $D_0(A, B)$  can be computed in combined linear time  $O(|A| \cdot |B|)$  from automata  $A$  and  $B$ , so that its size is in  $O(|A| \cdot |B|)$ . Furthermore, we can compute the least fixed point  $\text{lfp}(D_0(A, B))$  in combined linear time, too (Theorem 3).

Whether  $A, B \models \text{fail}_1$  can be tested in  $O(|A| \cdot |B|)$  is non-trivial though. To this purpose, we introduce a binary predicate  $\text{frb}(p, q)$  of *forbidden states*, which is equivalent to the implication  $\text{acc}(p, q) \rightarrow \text{fail}_1$ , i.e.:

$$A, B \models \text{frb}(p, q) \Leftrightarrow_{\text{df}} A, B \models \text{acc}(p, q) \Rightarrow \text{fail}_1.$$

Inclusion is thus violated if forbidden states are accessible. The following lemma is an immediate consequence of the definitions.

**Lemma 6.**  $A, B \models \text{frb}(p, q)$  iff there are  $p', q'$  such that one of the following two conditions holds:

1.  $A \models p @ p' \wedge A, B \models \text{acc}(p', q') \wedge B \not\models q @ q'$ , or
2.  $A \models p' @ p \wedge A, B \models \text{acc}(p', q') \wedge B \not\models q' @ q$ .

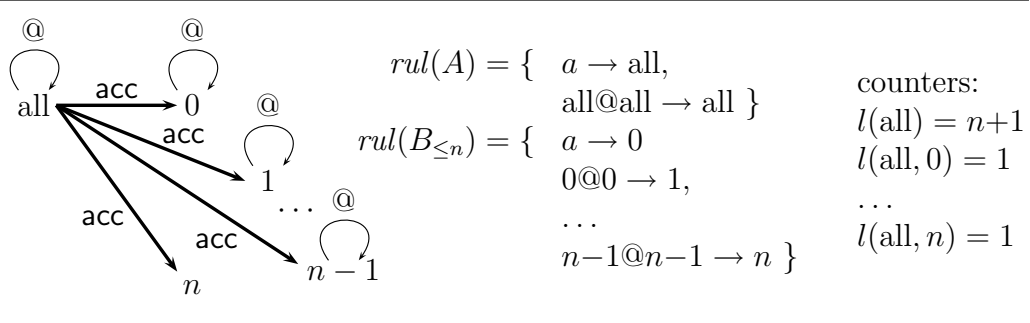


Figure 6: Rule  $(\text{frb}_{/2})$  can infer, for all  $0 \leq i \neq j \leq n$ , the clause  $\text{frb}(\text{all}, i) :- \text{acc}(\text{all}, j)$ .

Figure 5 extends  $D_0(A, B)$  to  $D_1(A, B)$  by clauses for  $\text{fail}_1$ . These are produced by transformation rules  $(\text{frb}_{/1})$  and  $(\text{frb}_{/2})$ , that are justified by Lemma 6. Transformation rule  $(\text{fail}_1)$  is witnessed by the definition of  $A, B \models \text{frb}(p, q)$ . It should be clear that  $A, B \models \text{fail}_1$  iff  $\text{fail}_1 \in \text{lfp}(D_1(A, B))$ .

**Proposition 7.** *Let  $A$  and  $B$  be stepwise tree automata for binary trees. If  $A$  is productive and  $B$  deterministic then:*

$$L(A) \subseteq L(B) \Leftrightarrow \text{lfp}(D_1(A, B)) \cap \{\text{fail}_0, \text{fail}_1, \text{fail}_2\} = \emptyset$$

*Proof.* From Proposition 5 and Lemma 6. □

However, the number of clauses produced by transformation rules  $(\text{frb}_{/1})$  and  $(\text{frb}_{/2})$  may sum up to  $O(|A| \cdot |\text{sta}(B)|^2)$  in the worst case. The overall size of the Datalog program  $D_1(A, B)$  is thus bounded by  $O(|A| \cdot |B|^2)$ . Programs of this size may arise, as shown by the example in Figure 6. Even though in this case  $O(|A| \cdot |B|) = O(n)$ , there are  $n^2$  clauses in  $D_1(A, B)$  produced by transformation rule  $(\text{frb}_{/2})$ .

By computing the least fixed point of  $D_1(A, B)$ , we can thus decide language inclusion  $L(A) \subseteq L(B)$  in time  $O(|A| \cdot |B|^2)$ . Unfortunately, this is not yet any better than the naive algorithm through complementation discussed in the introduction.

### 3.4. Inclusion Test in Time $O(|A| \cdot |B|)$

The problem with the clauses produced by  $(\text{frb}_{/1})$  and  $(\text{frb}_{/2})$  is the enumeration of clauses for forbidden states. This operation makes negative information of  $B$  explicit that one would like to leave implicit.

To see this, let  $\Sigma = \{a\}$  and let us consider the example in Figure 6. There, automaton  $A$  has a unique (final) state such that  $\text{sta}(A) = \text{fin}(A) = \{\text{all}\}$ . It recognizes all binary trees over  $\Sigma_{@}$  whose

paths to the left are of arbitrary length. They can be obtained as Curried encodings of unranked trees whose nodes have arbitrary many children. Automaton  $B_{\leq n}$  has states  $sta(B_{\leq n}) = fin(B_{\leq n}) = \{0, \dots, n\}$ . It recognizes the set of binary trees over  $\Sigma_{@}$  whose paths to the left are bounded by  $n$ . Thus, their Curried encodings cannot have more than  $n$  children. Note that the rules such that  $A \models p@p$ , resp.  $B_{\leq n} \models q@q$ , are depicted in Figure 6 by @-loops on states  $p \in sta(A)$ , resp.  $q \in sta(B_{\leq n})$ . Accessibility  $A, B_{\leq n} \models acc(all, j)$  holds for all  $0 \leq j \leq n$  and implies forbidden states  $A, B_{\leq n} \models frb(all, i)$  for all  $0 \leq i \leq n$ . This is inferred by quadratically many clauses  $frb(all, i) :- acc(all, j)$  that need to be avoided.

The idea is to count positive information in order to deduce how many times negative information can be inferred. Given a state  $p$ , we count the number of pairs  $(p', q')$  with  $A, B \models acc(p', q')$  and  $A \models p'@p$  or *vice versa*, and compare it with the number of such pairs that raise  $frb(p, q)$ :

$$\begin{aligned} l(p) &= \#\{(p', q') \mid A \models p'@p \wedge A, B \models acc(p', q')\} \\ &\quad + \#\{(p', q') \mid A \models p@p' \wedge A, B \models acc(p', q')\} \\ l(p, q) &= \#\{(p', q') \mid A \models p'@p \wedge A, B \models acc(p', q') \wedge B \models q'@q\} \\ &\quad + \#\{(p', q') \mid A \models p@p' \wedge A, B \models acc(p', q') \wedge B \models q@q'\} \end{aligned}$$

**Lemma 8.**  $A, B \models frb(p, q)$  iff  $l(p) > l(p, q)$ .

*Proof.* By definition,  $l(p) \geq l(p, q)$  for all  $p, q$ . We have  $l(p) > l(p, q)$  iff there are  $p', q'$  such that  $A \models p'@p \wedge A, B \models acc(p', q') \wedge B \not\models q'@q$  or symmetrically  $A \models p@p' \wedge A, B \models acc(p', q') \wedge B \not\models q@q'$ . By Lemma 6, this is equivalent to  $A, B \models frb(p, q)$ .  $\square$

It remains to see that we can compute the collection of numbers  $l(p)$  and  $l(p, q)$  for all  $p \in sta(A)$  and  $q \in sta(B)$  in time  $O(|A| \cdot |B|)$ . This can be done by the algorithm in Figure 7.

**Theorem 9.** Let  $A$  and  $B$  be tree automata over a ranked signature  $\Sigma$ , possibly with  $\epsilon$ -rules in  $A$ . If  $B$  is deterministic, then inclusion  $L(A) \subseteq L(B)$  can be decided in time  $O(|A| \cdot |B|)$  independently of the size of  $\Sigma$ .

*Proof.* We can assume that  $A$  and  $B$  are stepwise tree automata by Proposition 1. We first compute  $D_0(A, B)$  from  $A$  and  $B$  in combined linear time  $O(|A| \cdot |B|)$ , and then the least fixed point of this Datalog program in the same time. If  $lfp(D_0(A, B))$  contains  $fail_0$  or  $fail_2$ , then inclusion  $L(A) \subseteq L(B)$  does not hold. Otherwise, we compute all numbers  $l(p)$  and  $l(p, q)$  in time  $O(|A| \cdot |B|)$  by the algorithm in Figure 7, and test for all  $acc(p, q) \in lfp(D_0(A, B))$  whether  $A, B \models frb(p, q)$ . Inclusion  $L(A) \subseteq L(B)$  holds iff this test succeeds. It can be performed in time  $O(|A| \cdot |B|)$  by checking the values of the counters (Lemma 8).  $\square$

---

```

for all  $p \in sta(A)$  do  $l(p) := 0$ ;
  for all  $q \in sta(B)$  do  $l(p, q) := 0$ ;
for all  $p_1 @ p_2 \rightarrow p$  in  $rul(A)$  do
  for all  $q \in sta(B)$  do
    if  $acc(p_1, q) \in lfp(D_0(A, B))$  then increment  $l(p_2)$ ;
    if  $acc(p_2, q) \in lfp(D_0(A, B))$  then increment  $l(p_1)$ ;
  for all  $q_1 @ q_2 \rightarrow q$  in  $rul(B)$  do
    if  $acc(p_1, q_1) \in lfp(D_0(A, B))$  then increment  $l(p_2, q_2)$ ;
    if  $acc(p_2, q_2) \in lfp(D_0(A, B))$  then increment  $l(p_1, q_1)$ ;

```

---

Figure 7: Counting in  $O(|A| \cdot |B|)$ .

### 3.5. Efficient Algorithm

The previous algorithm has a satisfactory worst case complexity of  $O(|A| \cdot |B|)$ . In practice, however, it is non-optimal with respect to average time efficiency.

The first problem is that all pairs of rules in  $A$  and  $B$  are enumerated when computing the values of the counters. We now present a better algorithm, which inspects at most the accessible part of  $A \times B$ . The second problem is that  $(fail_1)$  is applied only after the fixed point computation. From now on, we envisage an algorithm (presented in full in Section 3.6) that detects inclusion failure as early as possible so that we do not have to complete the fixed point computation in such cases. These cases are very frequent in practice, as we will show experimentally (in Section 6), so that the gain in efficiency is considerable.

We introduce literals  $\mathit{frb}^c(p, Q)$  for states  $p \in sta(A)$  and state sets  $Q \subseteq sta(B)$  with the following semantics:

$$A, B \models \mathit{frb}^c(p, Q) \Leftrightarrow_{df} \forall q \in sta(B) \setminus Q. A, B \models \mathit{frb}(p, q)$$

In Figure 6, we have  $A, B \models \mathit{frb}^c(\mathit{all}, \{i\})$  for all  $0 \leq i \leq n$ . Thus, all literals  $\mathit{frb}(\mathit{all}, i)$  are implied by  $n$  literals  $\mathit{frb}^c(\mathit{all}, \{j\})$ . This multiplicity  $n$  is equal to  $l(\mathit{all}) - l(\mathit{all}, i)$ . Indeed, our objective is to compute the set of all literals satisfying  $A, B \models \mathit{frb}^c(p, Q)$  by a Datalog program and to infer the values of  $l(p) - l(p, q)$  thereby.

Note that two literals  $\mathit{frb}^c(p, Q)$  and  $\mathit{frb}^c(p, Q')$  are equal syntactically if and only if  $Q = Q'$ . In order to make this happen technically, we assume a fixed total order  $<$  on  $sta(B)$  in order to identify  $\mathit{frb}^c(p, Q)$  with the unique  $(n + 1)$ -ary literal  $\mathit{frb}^c(p, q_1, \dots, q_n)$  with  $Q = \{q_1, \dots, q_n\}$  and  $q_1 < \dots < q_n$ . Thereby, all results on ground Datalog programs continue to apply.



---


$$\begin{array}{l}
(\text{frb}_{/1}^c) \frac{A \models p_1 @ p_2 \quad q_2 \in \text{sta}(B)}{\text{frb}^c(p_1, Q_1^B(q_2)) :- \text{acc}(p_2, q_2).} \quad Q_1^B(q_2) = \{q_1 \mid B \models q_1 @ q_2\} \\
(\text{frb}_{/2}^c) \frac{A \models p_1 @ p_2 \quad q_1 \in \text{sta}(B)}{\text{frb}^c(p_2, Q_2^B(q_1)) :- \text{acc}(p_1, q_1).} \quad Q_2^B(q_1) = \{q_2 \mid B \models q_1 @ q_2\}
\end{array}$$


---

Figure 8: Grouping clauses from  $(\text{frb}_{/1})$  and  $(\text{frb}_{/2})$ .

In Figure 8, we present transformation rules  $(\text{frb}_{/1}^c)$  and  $(\text{frb}_{/2}^c)$ , which produce Datalog clauses inferring  $\text{frb}^c(p, Q)$  literals. They group many clauses produced by a transformation rule  $(\text{frb}_{/i})$ . Consider  $i = 2$ . The transformation rule assumes  $A \models p_1 @ p_2$  and a state  $q_1 \in \text{sta}(B)$ . It then computes the set  $Q_2^B(q_1)$  of all states  $q_2$  with  $B \models q_1 @ q_2$ , and produces the clause  $\text{frb}^c(p_2, Q_2^B(q_1)) :- \text{acc}(p_1, q_1)$ . This is correct, since if  $A, B \models \text{acc}(p_1, q_1)$ , then for all  $q_2 \notin Q_2^B(q_1)$ , we have  $A, B \models \text{frb}(p_2, q_2)$  and thus  $A, B \models \text{frb}^c(p_2, Q_2^B(q_1))$ . In Figure 6, for instance, transformation  $(\text{frb}_{/1}^c)$  produces for all  $0 \leq i \leq n$  the clauses  $\text{frb}^c(\text{all}, \{i\}) :- \text{acc}(\text{all}, i)$ . The overall size of these clauses is linear in  $n$ , no more quadratic!

Let  $D_2(A, B)$  be the ground Datalog program which extends  $D_0(A, B)$  by the clauses from  $(\text{frb}_{/1}^c)$  and  $(\text{frb}_{/2}^c)$ . This program remains incomplete, in that  $\text{frb}^c(p, Q)$  literals are never used in order to infer  $\text{fail}_1$ .

**Lemma 10.**  $D_2(A, B)$  can be computed in time  $O(|A| \cdot |B|)$  from  $A$  and  $B$ .

*Proof.* We have seen the result for  $D_0(A, B)$  already. The number of clauses produced by transformation rule  $(\text{frb}_{/1}^c)$  is in  $O(|A| \cdot |\text{sta}(B)|)$  but the size of each such clause is  $n + 1$  which in the worst case could be  $|\text{sta}(B)| + 1$ . Symmetrically for  $(\text{frb}_{/2}^c)$ . The overall size of all  $\text{frb}^c$  clauses, however, is bounded by the overall number of  $\text{acc}$  clauses produced at the same time, which in turn is bounded by  $O(|A| \cdot |B|)$ , too! To see this, we can rewrite the first rule of  $(\text{frb}_{/2}^c)$  as shown in Figure 9, such that the corresponding  $(\text{acc}_{/3})$  clauses are inferred simultaneously (and these don't overlap).

It remains to show how to compute  $D_2(A, B)$  in combined linear time. The following program produces all clauses from transformation rule  $(\text{frb}_{/2}^c)$ :

```

for all  $p_1 @ p_2 \rightarrow p \in \text{rul}(A)$  do
  for all  $q_1 \in \text{sta}(B)$  do
    compute  $Q = Q_2^B(q_1)$ ;
    collect  $\text{frb}^c(p_2, Q) :- \text{acc}(p_1, q_1)$ ;

```

The set  $Q$  can be computed in time  $O(|Q|)$  from a precomputed data structure that returns for a given state  $q_1$  all rules  $q_1 @ q_2 \rightarrow q$  in  $\text{rul}(B)$  in linear time depending on their number. The whole programs thus runs in time  $O(|D_2(A, B)|)$  which is in  $O(|A| \cdot |B|)$ .  $\square$

---


$$\left. \begin{array}{l} p_1 @ p_2 \rightarrow p \in \text{rul}(A) \\ q_1 \in \text{sta}(B) \end{array} \right\} \begin{array}{l} q_1 @ q_2^1 \rightarrow q^1 \in \text{rul}(B) \\ \vdots \\ q_1 @ q_2^n \rightarrow q^n \in \text{rul}(B) \end{array} \left. \vphantom{\begin{array}{l} p_1 @ p_2 \rightarrow p \in \text{rul}(A) \\ q_1 \in \text{sta}(B) \end{array}} \right\} \begin{array}{l} \text{all rules} \\ \text{for } q_1 \end{array}$$


---


$$\begin{array}{l} \text{frb}^c(p_2, \{q_2^1, \dots, q_2^n\}) :- \text{acc}(p_1, q_1). \\ \text{acc}(p, q^1) :- \text{acc}(p_1, q_1), \text{acc}(p_2, q_2^1). \\ \vdots \\ \text{acc}(p, q^n) :- \text{acc}(p_1, q_1), \text{acc}(p_2, q_2^n). \end{array}$$


---

Figure 9: Rewriting groups of  $(\text{frb}_{/2})$  clauses to  $(\text{frb}_{/2}^c)$  clauses.

The Datalog programs  $D_1(A, B)$  and  $D_2(A, B)$  have the same clauses for literals  $\text{acc}(p, q)$ ,  $\text{fail}_0$ , and  $\text{fail}_2$ , so their least fixed points coincide for these. In particular, we can decide  $A, B \models \text{fail}_0 \vee \text{fail}_2$  in time  $O(|A| \cdot |B|)$  by testing membership of  $\text{fail}_0$  and  $\text{fail}_2$  in  $\text{lfp}(D_2(A, B))$ . It remains to relate both programs with respect to forbidden states and  $\text{fail}_1$ .

**Lemma 11.** *A literal  $\text{frb}(p, q)$  belongs to  $\text{lfp}(D_1(A, B))$  if and only if there exists a set  $Q \subseteq \text{sta}(B)$  not containing  $q$  such that  $\text{frb}^c(p, Q) \in \text{lfp}(D_2(A, B))$ .*

*Proof.* Suppose that  $\text{frb}(p_1, q_1) \in \text{lfp}(D_1(A, B))$ . The previous literal has been added by a clause produced by  $(\text{frb}_{/1})$  or  $(\text{frb}_{/2})$ . By symmetry it is sufficient to consider the first case only. The contributing clause of  $D_1(A, B)$  must be of the form  $\text{frb}(p_1, q_1) :- \text{acc}(p_2, q_2)$ .  $(\text{frb}_{/1})$  assumes  $A \models p_1 @ p_2$  and  $B \not\models q_1 @ q_2$ , so that  $q_1 \notin Q_1(q_2)$ .  $(\text{frb}_{/1}^c)$  produces the clause  $\text{frb}^c(p_1, Q_1^B(q_2)) :- \text{acc}(p_2, q_2)$  in  $D_2(A, B)$ . Since  $\text{acc}(p_2, q_2) \in \text{lfp}(D_1(A, B))$ , we equally have  $\text{acc}(p_2, q_2) \in \text{lfp}(D_2(A, B))$ . Hence, the above clause of  $D_2(A, B)$  is applicable and adds  $\text{frb}^c(p_1, Q_1^B(q_2))$  to  $\text{lfp}(D_2(A, B))$ . The inverse argument is similar.  $\square$

**Lemma 12.** *For  $D = D_2(A, B)$ ,  $p \in \text{sta}(A)$  and  $q \in \text{sta}(B)$ :*

$$\begin{aligned} l(p) &= \sum_{Q \subseteq \text{sta}(B)} \text{lfp}^\#(D)(\text{frb}^c(p, Q)) \\ l(p, q) &= \sum_{Q \subseteq \text{sta}(B), q \in Q} \text{lfp}^\#(D)(\text{frb}^c(p, Q)) \end{aligned}$$

*Proof.* This follows from the definitions, as we elaborate in the first case:

$$\begin{aligned} l(p) &= \#\{(p', q') \mid A \models p' @ p \wedge A, B \models \text{acc}(p', q')\} \\ &+ \#\{(p', q') \mid A \models p @ p' \wedge A, B \models \text{acc}(p', q')\} \\ &= \#\{(p', q') \mid \text{frb}^c(p, Q_2^B(q')) :- \text{acc}(p', q') \in D \wedge \text{acc}(p', q') \in \text{lfp}(D)\} \\ &+ \#\{(p', q') \mid \text{frb}^c(p, Q_1^B(q')) :- \text{acc}(p', q') \in D \wedge \text{acc}(p', q') \in \text{lfp}(D)\} \\ &= \sum_{Q \subseteq \text{sta}(B)} \text{lfp}^\#(D)(\text{frb}^c(p, Q)) \end{aligned} \quad \square$$

We can thus compute all numbers  $l(p)$  and  $l(p, q)$  in linear time depending on the size of  $lfp(D_2(A, B))$  by Corollary 4.

Inclusion  $L(A) \subseteq L(B)$  can be tested as before, except that the numbers  $l(p)$  and  $l(p, q)$  can now be computed from  $lfp(D_2(A, B))$  more efficiently. This requires computing the accessible part of  $A \times B$  only, by lazily creating only the needed clauses from  $(acc_{/3})$  as usual. The application of all rules  $(frb_{/1}^c)$  and  $(frb_{/2}^c)$  can be done in time  $O(|A| \cdot |sta(B)| + |rul(B)|)$ , which may be much smaller than  $O(|A \times B|)$ , too.

### 3.6. Early Failure Detection

We next tackle the problem that  $(fail_1)$  is tested only after fixed point computation at the end, by the following loop:

```
for all  $acc(p, q) \in lfp(D_2(A, B))$  do
  if  $l(p) > l(p, q)$  then return false;
  otherwise return true;
```

Instead, we approach on the fly checking for  $(fail_1)$  as follows. Once a literal  $acc(p, q)$  is inferred during the computation of the fixed point of  $D_2(A, B)$ , we check in constant time whether the current values of the counters satisfy  $l(p) > l(p, q)$ , i.e., whether  $frb(p, q)$  is implied by some literal  $frb^c(p, Q)$  inferred before with  $q \notin Q$ . This requires updating the counters on the fly, but this is not difficult if we update them with priority.

The main difficulty arises when deriving  $frb^c(p, Q)$  only after some  $acc(p, q)$ , since we cannot check for all  $q \in sta(B) \setminus Q$  whether  $acc(p, q)$  has been inferred before without enumerating the complement of  $Q$ . It turns out fortunately that all tests for  $(fail_1)$  come for free and on the fly (without any testing after fixed point computation) if we impose the following *priority* discipline. We assume that literals of the form  $acc(p, q)$  are always inferred with the lowest priority, i.e., whenever other literals can be inferred at the same time, these will be inferred before.

Our *on-the-fly algorithm* thus computes the least fixed point of  $D_2(A, F)$  with the above priorities. The counters  $l(p)$  and  $l(p, q)$  are always updated immediately. Whenever a literal  $acc(p, q)$  is inferred, the counters are tested for  $l(p) > l(p, q)$ . If this test succeeds, the algorithm returns false, otherwise it continues and returns true at the very end.

**Lemma 13.** *The on-the-fly algorithm correctly detects  $(fail_1)$  if a literal  $acc(p_1, q_1)$  is inferred before some literal  $frb^c(p_1, Q_1)$  with  $q_1 \notin Q_1$ .*

*Proof.* This situation is depicted in Figure 10. Literal  $frb^c(p_1, Q_1)$  originates from a clause produced by rule  $(frb^c)$  and some literal  $acc(p_2, q_2)$  added earlier:

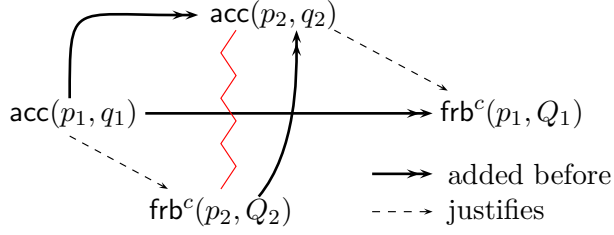


Figure 10: Early failure detection:  $\text{frb}^c(p_2, Q_2)$  in  $\text{lf}p(D_2(A, B))$  before  $\text{acc}(p_2, q_2)$ .

$$\frac{p_1 @ p_2 \rightarrow p \in \text{rul}(A) \quad Q_1 = Q_1^B(q_2)}{\text{frb}^c(p_1, Q_1) :- \text{acc}(p_2, q_2)}.$$

We show by contradiction that  $\text{acc}(p_1, q_1)$  has got added before  $\text{acc}(p_2, q_2)$ . Otherwise,  $\text{acc}(p_2, q_2)$  has been added before  $\text{acc}(p_1, q_1)$ , so that due to our priority assumption,  $\text{frb}^c(p_1, Q_1)$  has been added before  $\text{acc}(p_1, q_1)$  which contradicts the hypothesis. Having  $\text{acc}(p_1, q_1)$  in the fixed point permits to apply the following clause of ( $\text{frb}^c$ ):

$$\frac{p_1 @ p_2 \rightarrow p \in \text{rul}(A) \quad Q_2 = Q_2^B(q_1)}{\text{frb}^c(p_2, Q_2) :- \text{acc}(p_1, q_1)}.$$

Note that  $q_1 \in Q_1^B(q_2)$  iff  $q_2 \in Q_2^B(q_1)$ . Thus  $q_2 \notin Q_2$  since  $q_1 \notin Q_1$ . Consequently,  $\text{acc}(p_2, q_2), \text{frb}^c(p_2, Q_2) \in \text{lf}p(D_2(A, B))$  raises ( $\text{fail}_1$ ), and this is correctly detected by the modified algorithm, since  $\text{frb}^c(p_2, Q_2)$  is inferred before  $\text{acc}(p_2, q_2)$ .  $\square$

### 3.7. Incrementality

Incrementality appears to be critical for efficiency in our experiments. In our prime application to schema-guided query induction [9], for instance, we use incremental addition of  $\epsilon$ -rules to the automaton  $A$  on the left. These model state merging operations  $p_1 = p_2$  during automata induction as  $p_1 \xrightarrow{\epsilon} p_2$  and  $p_2 \xrightarrow{\epsilon} p_1$ .

Fixed points of Datalog programs can be computed incrementally with respect to adding new clauses. Priorities, however, may raise trouble here. It would not be correct to add clauses later on, that should have been applied with priority before. In this case, one would have to redo some work.

Rules of automata  $A$  or  $B$  are transformed to clauses of  $D_2(A, B)$ . The incremental addition of  $\epsilon$ -rules to  $A$  is harmless. They are transformed by ( $\text{acc}_{/2}$ ) to clauses with the lowest priority. All previous clauses remain valid (in contrast to adding rules  $q_1 @ q_2 \rightarrow q$  to  $B$  which changes  $Q_2^B(q_1)$ ). For these two reasons, we do not have to redo any work when adding  $\epsilon$ -rules to  $A$  later on. Of course, incrementality assumes early failure detection.

## 4. Factorized Tree Automata

We next relax the determinism assumption on  $B$  in a controlled manner, that will be crucial to deal with DTDs. This leads us to introduce the notion of deterministic factorized tree automata, and to check inclusion for them. Inclusion in deterministic factorized automata is exactly what we need for inclusion in deterministic DTDs in Section 5.

### 4.1. Deterministic Factorized Tree Automata

We replace  $B$  by deterministic factorized automata  $F$ , which we now introduce. These are stepwise tree automata with  $\epsilon$ -rules for ranked trees, that represent deterministic stepwise tree automata in a more compact manner.

**Definition 1.** A factorized tree automaton  $F$  over a stepwise signature  $\Sigma$  consists of a stepwise tree automaton with  $\epsilon$ -rules and a partition  $sta(F) = sta_1(F) \uplus sta_2(F)$  such that for all  $q_1 @ q_2 \rightarrow q$  in  $rul(F)$  we have  $q_1 \in sta_1(F)$  and  $q_2 \in sta_2(F)$ .

We say that  $q$  is of sort  $i$  in  $F$  if  $q \in sta_i(F)$ . The sort determines which states may be used in the  $i$ th position of the binary symbol  $@$  in rules of  $F$ .

Every factorized automaton  $F$  defines a tree automaton  $b(F)$  without  $\epsilon$ -rules that recognizes the same language. Both automata have the same signature and states; the rules of  $b(F)$  are inferred as follows from those of  $F$ :

$$(E_1) \frac{a \rightarrow q \in rul(F)}{a \rightarrow q \in rul(b(F))} \quad (E_2) \frac{q_1 \xrightarrow{\epsilon}_F^* r_1 \quad q_2 \xrightarrow{\epsilon}_F^* r_2 \quad r_1 @ r_2 \rightarrow q \in rul(F)}{q_1 @ q_2 \rightarrow q \in rul(b(F))}$$

We set  $fin(b(F)) = \{q \mid q \xrightarrow{\epsilon}_F^* r, r \in fin(F)\}$ . Note that the size of  $b(F)$  may be  $O(|rul(F)| \cdot |sta(F)|^2)$ , which is cubic in that of  $F$  in the worst case. Besides their succinctness, the truly interesting bit about factorized tree automata is their notion of determinism.

A collection of examples for factorized tree automata  $(F_{\leq n})_n$  is given in Figure 11. The set of constants of  $F_{\leq n}$  is  $\Sigma = \{a, b, c, d, e, f\}$ . Automaton  $F_{\leq n}$  recognizes all binary trees over  $\Sigma_{@}$  whose paths to the left are always of length at most  $n$ . These can be obtained as Curried encodings of unranked trees whose nodes have at most  $n$  children. The states of sort 1 of  $F_{\leq n}$  are  $\{0, \dots, n\}$ . For every node, they count the length of the path to the left-most leaf. The single state of sort 2 is  $ok$ . It can be assigned to all nodes rooting subtrees in the language of  $F_{\leq n}$ . Automaton  $F_{\leq n}$  has rules  $a, b, c, d, e, f \rightarrow 0$ , rules  $i-1 @ ok \rightarrow i$  for all  $1 \leq i \leq n$  and  $i \xrightarrow{\epsilon} ok$  for all  $0 \leq i \leq n$ . The

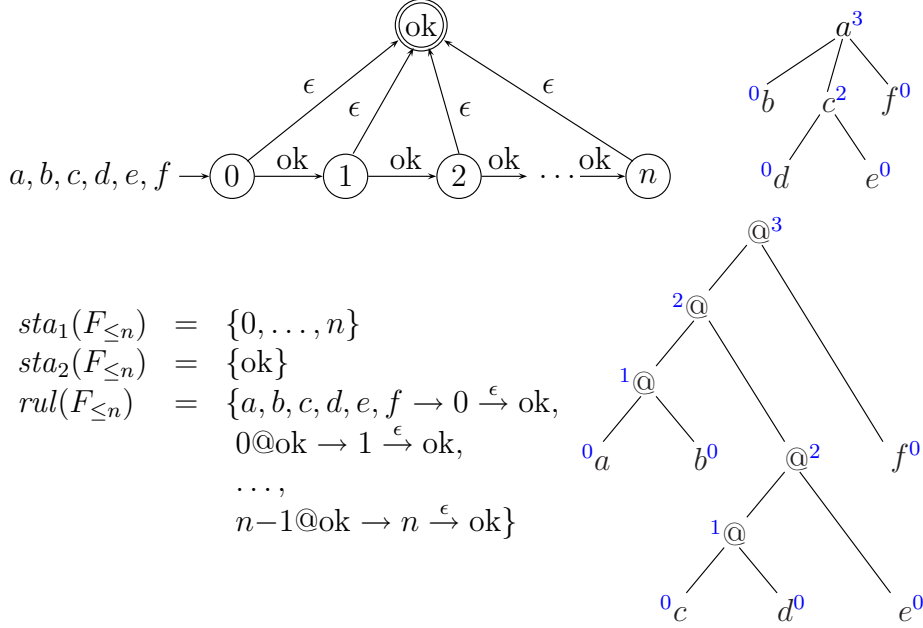


Figure 11: On the left, we define for all  $n$  a deterministic factorized tree automaton  $F_{\leq n}$ . On the right, the lower tree is the Curried encoding of the upper. We annotate all nodes by the unique state assigned by the evaluator of the deterministic tree automaton  $b(F_{\leq n})$  where  $n \geq 3$ .

size of  $F_{\leq n}$  is thus in  $O(n)$ . The corresponding tree automaton  $b(F_{\leq n})$  is of size  $O(n^2)$ , since it has rules  $a, b, c, d, e, f \rightarrow 0$  and  $i-1@j \rightarrow i$  for all  $1 \leq i \leq n$  and  $0 \leq j \leq n$ . Note that  $b(F_{\leq n})$  is the unique state-minimal deterministic automaton recognizing the language of  $F_{\leq n}$ . The sizes  $|F_{\leq n}|$  are asymptotically smaller by a factor of  $n$  than  $|b(F_{\leq n})|$ . Nevertheless, all  $F_{\leq n}$  are deterministic in the following sense.

**Definition 2.** A factorized tree automaton  $F$  is (bottom-up) deterministic if:

$\mathbf{d}_0$ : the  $\epsilon$ -free part of  $F$  is (bottom-up) deterministic;

$\mathbf{d}_1$ : for all  $q \in sta(F)$  and sorts  $i \in \{1, 2\}$ , there is at most one state  $r$  of sort  $i$  such that  $q \xrightarrow{\epsilon^*}_F r$ .

Non-redundant  $\epsilon$ -rules must change the sort: if  $q \xrightarrow{\epsilon}_F r$  for two states of the same sort, then  $r = q$  by  $\mathbf{d}_1$ , and  $q \xrightarrow{\epsilon^*}_F q$ . A similar argument shows that all proper chains of  $\epsilon$ -rules are redundant so that  $\xrightarrow{\epsilon^*}_F$  is equal to  $\xrightarrow{\epsilon^{\leq 1}}_F$ .

Another consequence of determinism is that the size of deterministic  $b(F)$  is at most quadratic in the size of deterministic  $F$ , since  $b(F)$  cannot have more than  $|sta(F)|^2$  many binary transitions.

**Proposition 14.** *The tree automaton  $b(F)$  is deterministic for all deterministic factorized tree automata  $F$ .*

*Proof.* Let  $B = b(F)$  which by construction is free of  $\epsilon$ -rules. For every constant  $a \in \Sigma$ , the uniqueness of  $q$  such that  $a \rightarrow q \in rul(B)$  follows from  $\mathbf{d}_0$ . For every  $q_1 @ q_2 \rightarrow q$  in  $rul(B)$  we have to show that  $q$  is uniquely determined by  $q_1$  and  $q_2$ . By  $\mathbf{d}_1$  there is at most one state  $r_1$  of sort 1 such that  $q_1 \xrightarrow{\epsilon}^* r_1$  and at most one state  $r_2$  of sort 2 such that  $q_2 \xrightarrow{\epsilon}^* r_2$ . Condition  $\mathbf{d}_0$  implies that there exists at most one state  $q$  such that  $r_1 @ r_2 \rightarrow q \in rul(F)$ .  $\square$

Conversely, every deterministic stepwise tree automaton  $B$  can be converted in time  $O(|B|)$  into a deterministic factorized tree automaton  $F$ , such that  $b(F)$  is equal to  $B$  modulo state renaming. The states of  $F$  are  $sta_i(F) = sta(B) \times \{i\}$  for  $i \in \{1, 2\}$ . Rules are transformed as follows:

$$\frac{q_1 @ q_2 \rightarrow q \in rul(B)}{(q_1, 1) @ (q_2, 2) \rightarrow (q, 1) \in rul(F)} \quad \frac{a \rightarrow q \in rul(B)}{a \rightarrow (q, 1) \in rul(F)}$$

$$(q, 1) \xrightarrow{\epsilon} (q, 2) \in rul(F) \quad (q, 1) \xrightarrow{\epsilon} (q, 2) \in rul(F)$$

#### 4.2. Testing Validity of $\mathbf{fail}_0$ and $\mathbf{fail}_2$

Given an automaton  $A$  and a deterministic factorized automaton  $F$ , we first characterize  $A, b(F) \models \mathbf{fail}_0$  and  $A, b(F) \models \mathbf{fail}_2$  in terms of  $A$  and  $F$ . This must be done without computing  $b(F)$ , since its size may be in  $O(|sta(F)|^2)$ .

**Lemma 15.** *Let  $B = b(F)$ .*

- (1)  $A, B \models \mathbf{fail}_0$  iff  $\exists a \rightarrow p \in rul(A) \wedge \nexists a \rightarrow q \in rul(F)$
- (2)  $A, B \models \mathbf{fail}_2$  iff  $\exists p \in fin(A) \exists q \in sta(F). A, B \models \mathbf{acc}(p, q) \wedge \forall r \in sta(F). q \xrightarrow{\epsilon}^{\leq 1} r \Rightarrow r \notin fin(F)$

*Proof.* The first statement follows from construction rule  $(E_1)$  of  $b(F)$ . For the second, note that  $q \notin fin(B)$  iff  $\forall r \in sta(F). q \xrightarrow{\epsilon}^{\leq 1} r \Rightarrow r \notin fin(F)$ . Note that this universal quantifier is harmless, since for every  $q$  there is at most one  $r$  with  $q \xrightarrow{\epsilon}^{\leq 1} r$ . We can now conclude straightforwardly:

$$A, B \models \mathbf{fail}_2$$

$$\Leftrightarrow \exists p \in fin(A) \exists q \notin fin(B). A, B \models \mathbf{acc}(p, q)$$

$$\Leftrightarrow \exists p \in fin(A) \exists q \in sta(F). A, B \models \mathbf{acc}(p, q) \wedge \forall r. q \xrightarrow{\epsilon}^{\leq 1} r \Rightarrow r \notin fin(F) \quad \square$$

There is a subtle difference between accessibility in  $F$  and  $b(F)$ : accessibility in  $b(F)$  implies accessibility in  $F$ , but not *vice versa*. For instance, in

---


$$\begin{array}{l}
(\text{acc}/_{1a}) \frac{a \rightarrow p \in \text{rul}(A) \quad a \rightarrow q \in \text{rul}(F)}{\text{acc}(p, q)}. \\
(\text{acc}/_{2a}) \frac{p' \xrightarrow{\epsilon} p \in \text{rul}(A) \quad q \in \text{sta}(F)}{\text{acc}(p, q) :- \text{acc}(p', q)}. \\
(\text{acc}/_{3a}) \frac{p_1 @ p_2 \rightarrow p \in \text{rul}(A) \quad q_1 @ q_2 \rightarrow q \in \text{rul}(F)}{\text{acc}(p, q) :- \text{f.acc}(p_1, q_1), \text{f.acc}(p_2, q_2)}. \\
(\text{f.acc}) \frac{p \in \text{sta}(A) \quad q \xrightarrow{\epsilon \leq 1} r}{\text{f.acc}(p, r) :- \text{acc}(p, q)}. \\
(\text{fail}_{0/a}) \frac{a \rightarrow p \in \text{rul}(A) \quad \nexists q \in \text{sta}(F). a \rightarrow q \in \text{rul}(F)}{\text{fail}_0}. \\
(\text{fail}_{2/a}) \frac{p \in \text{fin}(A) \quad \forall r \in \text{sta}(F). q \xrightarrow{\epsilon \leq 1} r \Rightarrow r \notin \text{fin}(F)}{\text{fail}_2 :- \text{acc}(p, q)}.
\end{array}$$


---

Figure 12:  $D_0(A, F)$  testing  $A, B \models \text{fail}_0$  and  $A, B \models \text{fail}_2$  where  $B = b(F)$ .

Figure 11, state `ok` is accessible in  $F_{\leq n}$  but not in  $b(F_{\leq n})$ . This illustrates a detail of the construction of  $b(F)$ , which is essential for the preservation of determinism (Proposition 14). This difference is inherited to accessibility in  $A \times F$  and  $A \times b(F)$ . In order to avoid ambiguities, we write  $A, F \models \text{f.acc}(p, q)$  if  $(p, q)$  is accessible in  $A \times F$ . Thus, the following implication holds but not its converse:

$$A, b(F) \models \text{acc}(p, q) \Rightarrow A, F \models \text{f.acc}(p, q)$$

We define a ground Datalog program  $D_0(A, F)$  in Figure 12 in order to compute all valid `acc` and `f.acc` literals, i.e., all literals with  $A, b(F) \models \text{acc}(p, q)$  and  $A, F \models \text{f.acc}(p, q)$ . Furthermore, program  $D_0(A, F)$  provides rules  $(\text{fail}_{0/a})$  and  $(\text{fail}_{2/a})$ , which infer literals `fail0` and `fail2` respectively, according to Lemma 15.

It remains to verify that the program  $D_0(A, F)$  does indeed infer all valid `acc` and `f.acc` literals. This is shown by the following Lemma.

**Lemma 16.** *Let  $B = b(F)$  and  $L = \text{lf}p(D_0(A, F))$ .*

1.  $A, B \models \text{acc}(p, q)$  iff  $\text{acc}(p, q) \in L$ .
2.  $A, F \models \text{f.acc}(p, q)$  iff  $\text{f.acc}(p, q) \in L$ .
3.  $A, B \models \text{fail}_0$  iff  $\text{fail}_0 \in L$ .
4.  $A, B \models \text{fail}_2$  iff  $\text{fail}_2 \in L$ .



*Proof.* The 4 implications from the right to the left can be shown by simultaneous induction of the definition of the least fixed point. This is technical but straightforward. For the implications from the left to the right, we proceed as follows.

1. We show that  $(p, q) \in eval_{A \times B}(t)$  implies  $\mathbf{acc}(p, q) \in L$  by induction on the structure of  $t$ . Here we need construction rule (E<sub>2</sub>) of  $b(F)$ .
2. We show that  $(p, q) \in eval_{A \times F}(t)$  implies  $\mathbf{f.acc}(p, q) \in L$  by induction on the structure of  $t$ .
3. From Lemma 15 and transformation rule (fail<sub>0/a</sub>).
4. From Lemma 15, transformation rule (fail<sub>2/a</sub>), and part (1) above.  $\square$

**Lemma 17.**  $D_0(A, F)$  can be computed in time  $O(|A| \cdot |F|)$  from  $A$  and  $F$ .

The proof is obvious. One can even compute the least fixed point of  $D_0(A, F)$  such that only the productive part of  $A \times F$  has to be inspected.

**Proposition 18.** We can test in time  $O(|A| \cdot |F|)$  whether  $A, b(F) \models \mathbf{fail}_0$  and  $A, b(F) \models \mathbf{fail}_2$ .

*Proof.* By Lemma 17 it is sufficient to compute the least fixed point of  $D_0(A, F)$  and to verify whether it contains  $\mathbf{fail}_0$ , resp.  $\mathbf{fail}_2$ . This can be done in time  $O(|A| \cdot |F|)$  by Lemma 17, even such that only the productive part of  $A \times F$  is inspected.  $\square$

### 4.3. Testing Validity of $\mathbf{fail}_1$

It remains to characterize  $A, b(F) \models \mathbf{fail}_1$  in terms of  $A$  and  $F$ . Our solution in Lemma 19 will be technically intricate.

We need literals for  $A, F$  which are  $\mathbf{f.frb}_1(p, q)$ ,  $\mathbf{f.frb}_2(p, q)$ , and  $\mathbf{frb}(p, q)$ , whose semantics is summarized in Figure 13. We start with  $\mathbf{f.frb}_1(p, q)$  literals:

$$A, F \models \mathbf{f.frb}_1(p_1, q_1) \Leftrightarrow_{df} \begin{cases} q_1 \in sta_1(F) \wedge \exists p_2, q_2. \\ A, F \models \mathbf{f.acc}(p_2, q_2) \wedge A \models p_1 @ p_2 \wedge F \not\models q_1 @ q_2 \end{cases}$$

Note that  $A, F \models \mathbf{f.frb}_1(p, q)$  does not always imply  $A, b(F) \models \mathbf{frb}(p, q)$ , since we do not require  $q_2 \in sta_2(F)$  in the above definition. The definition of  $A, F \models \mathbf{f.frb}_2(p, q)$  is symmetric. The third predicate has the following meaning, where  $R_i^F = \{q \mid \exists r \in sta_i(F). q \xrightarrow{F}^{\leq 1} r\}$  for  $i \in \{1, 2\}$ .

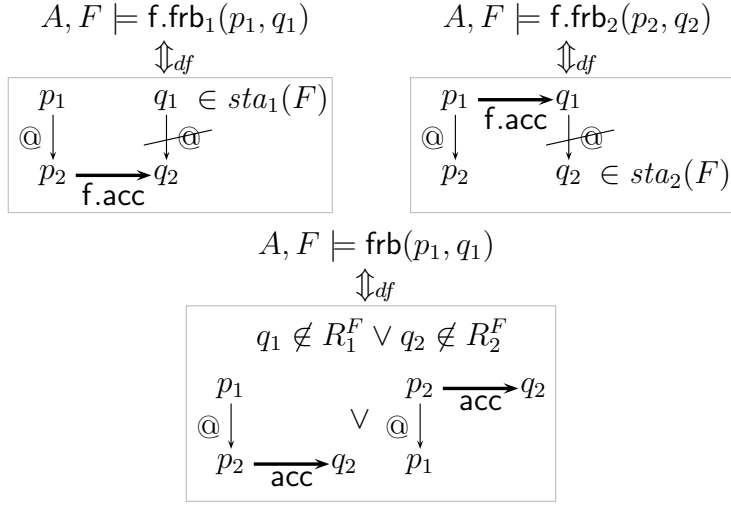


Figure 13: Semantics of predicates  $\text{f.frb}_1$ ,  $\text{f.frb}_2$  and  $\text{frb}$  for factorized tree automata.

$$A, F \models \text{frb}(p_1, q_1) \Leftrightarrow_{df} \begin{cases} \exists q_2. (q_1 \notin R_1^F \vee q_2 \notin R_2^F) \wedge \exists p_2. \\ A \models (p_1 @ p_2 \vee p_2 @ p_1) \wedge A, F \models \text{acc}(p_2, q_2) \end{cases}$$

As the proof of Lemma 19 will show, it holds that  $A, F \models \text{frb}(p, q)$  implies  $A, b(F) \models \text{frb}(p, q)$ , but not *vice versa*.

**Lemma 19.**  $A, b(F) \models \text{frb}(p, q)$  iff one of the following properties holds:

1. there exists  $i \in \{1, 2\}$  such that  $A, F \models \text{f.frb}_i(p, r)$ , where  $r$  is the unique state of sort  $i$  with  $q \xrightarrow{F}^{\leq 1} r$ , or
2.  $A, F \models \text{frb}(p, q)$ .

*Proof.* Let  $B = b(F)$ .

For the implication from the left to the right, we assume  $A, B \models \text{frb}(p_1, q_1)$ . By definition, there is a literal satisfying  $A, B \models \text{acc}(p_2, q_2)$  such that (a)  $A \models p_1 @ p_2$  and  $B \not\models q_1 @ q_2$ , or (b)  $A \models p_2 @ p_1$  and  $B \not\models q_2 @ q_1$ . By symmetry, it is sufficient to consider case (a). Part (1) of Lemma 16 shows that  $A, B \models \text{acc}(p_2, q_2)$  implies  $A, F \models \text{acc}(p_2, q_2)$ . We distinguish two exhaustive cases:

1. Case  $q_1 \in R_1^F \wedge q_2 \in R_2^F$ . There exists a unique state  $r_1 \in \text{sta}_1(F)$ , resp.  $r_2 \in \text{sta}_2(F)$ , such that  $q_1 \xrightarrow{F}^{\leq 1} r_1$ , resp.  $q_2 \xrightarrow{F}^{\leq 1} r_2$ . In this situation,  $B \not\models q_1 @ q_2$  is equivalent to  $F \not\models r_1 @ r_2$ . From  $A, B \models \text{acc}(p_2, q_2)$ , it follows that  $A, F \models \text{f.acc}(p_2, r_2)$  and hence,  $A, F \models \text{f.frb}_1(p_1, r_1)$ .

---


$$\begin{array}{l}
(\text{acc}/4) \frac{p \in \text{sta}(A) \quad q \in \text{sta}(F)}{\text{acc}(p, \_) :- \text{acc}(p, q)}. \\
(\text{frb}/1a) \frac{A \models p_1 @ p_2 \quad q_1 \notin R_1^F}{\text{frb}(p_1, q_1) :- \text{acc}(p_2, \_)}. \\
(\text{frb}/2a) \frac{A \models p_1 @ p_2 \quad q_2 \notin R_2^F}{\text{frb}(p_2, q_2) :- \text{acc}(p_1, \_)}. \\
(\text{fail}_1/a) \frac{p \in \text{sta}(A) \quad q \in \text{sta}(F)}{\text{fail}_1 :- \text{frb}(p, q), \text{acc}(p, q)}.
\end{array}$$


---

Figure 14:  $D_1(A, F)$  extends  $D_0(A, F)$  for checking  $\text{fail}_1$  raised by  $A, F \models \text{frb}(p, q)$ .

2. Case  $q_1 \notin R_1^F \vee q_2 \notin R_2^F$ . By definition, this implies  $A, F \models \text{frb}(p_1, q_1)$ .

For the other direction, we have to consider the two cases.

1. By symmetry, we can assume  $i = 1$ . We thus assume that the unique  $r_1 \in \text{sta}_1(F)$  with  $q_1 \xrightarrow{\epsilon \leq 1}_F r_1$  satisfies  $A, F \models \text{f.frb}_1(p_1, r_1)$ . By definition, there exist  $p_2$  and  $r_2$  such that  $A \models p_1 @ p_2$  and  $A, F \models \text{f.acc}(p_2, r_2)$ . By parts (2) and (1) of Lemma 16 there exists  $q_2$  such that  $A, B \models \text{acc}(p_2, q_2)$  and  $q_2 \xrightarrow{\epsilon \leq 1}_F r_2$ . In this situation,  $F \not\models r_1 @ r_2$  is equivalent to  $B \not\models q_1 @ q_2$ . Hence,  $A, B \models \text{frb}(p_1, q_1)$ .
2. We assume  $A, F \models \text{frb}(p_1, q_1)$  and show that  $A, B \models \text{frb}(p_1, q_1)$ . By definition, there exist  $q_2$  such that  $q_1 \notin R_1^F \vee q_2 \notin R_2^F$  and  $p_2$  such that  $A \models p_1 @ p_2 \vee p_2 @ p_1$  and  $A, B \models \text{acc}(p_2, q_2)$ . By symmetry, we can assume that  $A \models p_1 @ p_2$ . From  $q_1 \notin R_1^F \vee q_2 \notin R_2^F$ , it follows that  $B \not\models q_1 @ q_2$  and hence,  $A, B \models \text{frb}(p_1, q_1)$ .  $\square$

Our next goal is to test  $A, b(F) \models \text{fail}_1$ , when raised by  $A, F \models \text{frb}(p, q)$  and  $A, b(F) \models \text{acc}(p, q)$ , in time  $O(|A| \cdot |F|)$ . A naive Datalog program of size  $O(|A| \cdot |\text{sta}(F)|^2)$  is easy to deduce from the definition of  $A, F \models \text{frb}(p, q)$ . The less naive Datalog program  $D_1(A, F)$  in Figure 14 extends  $D_0(A, F)$ , in order to solve this task in time  $O(|A| \cdot |F|)$ . In order to avoid the quadratic factor, it relies on new literals  $\text{acc}(p, \_)$ , which we define to be equivalent to  $\exists q. \text{acc}(p, q)$ :

$$A, b(F) \models \text{acc}(p, \_) \Leftrightarrow_{df} A, F \models \exists q. \text{acc}(p, q)$$

All valid literals of type  $\text{acc}(p, \_)$  are computed by  $D_1(A, F)$  by clauses from  $(\text{acc}/4)$  and  $D_0(A, F)$ . The remaining clauses from  $D_1(A, F)$  check whether  $\text{fail}_1$  is raised by valid  $\text{frb}$  literals.

**Lemma 20.**  $A, F \models \exists p \exists q. \text{frb}(p, q) \wedge \text{acc}(p, q)$  iff  $\text{fail}_1 \in \text{lfp}(D_1(A, F))$ .

*Proof.* The soundness (“ $\Leftarrow$ ”) of the rules is obvious. It remains to show their completeness (“ $\Rightarrow$ ”). We assume  $A, F \models \text{frb}(p_1, q_1) \wedge \text{acc}(p_1, q_1)$ . By definition of  $A, F \models \text{frb}(p_1, q_1)$ , this holds in situations where  $A \models p_1 @ p_2$ ,  $A, F \models \text{acc}(p_2, q_2)$  and  $q_1 \notin R_1^F \vee q_2 \notin R_2^F$ . For symmetry, it is sufficient to consider the case  $q_1 \notin R_1^F$ . Let  $L = \text{lfp}(D_1(A, F))$ . Part (1) of Lemma 16 shows  $\text{acc}(p_1, q_1) \in L$  and hence  $\text{acc}(p_1, \_ ) \in L$  by  $(\text{acc}/_A)$ . From this, it can be deduced  $\text{frb}(p_2, q_2) \in L$  by  $(\text{frb}/_{1a})$ , so that  $\text{fail}_1 \in L$  by  $(\text{fail}_1/a)$ . Note that  $\text{frb}(p_1, q_1) \in L$  is possible, even though  $A, F \models \text{frb}(p_1, q_1)$ .  $\square$

**Lemma 21.**  $D_1(A, F)$  can be computed in time  $O(|A| \cdot |F|)$  from  $A$  and  $F$ .

*Proof.* All clauses depend on an element of  $A$  and an element of  $F$  only.  $\square$

Our next objective is to test  $A, F \models \text{f.frb}_i(p, q)$  for all  $p, q$  in time  $O(|A| \cdot |F|)$ . We consider  $i = 1$  only, for the sake of symmetry. Analogically to the case without factorization, we define the following counters:

$$\begin{aligned} l_1(p) &= \#\{(p', q') \mid A \models p' @ p \wedge A, F \models \text{f.acc}(p', q')\} \\ l_1(p, q) &= \#\{(p', q') \mid A \models p' @ p \wedge A, F \models \text{f.acc}(p', q') \wedge F \models q' @ q\} \end{aligned}$$

**Lemma 22.** For all  $q \in \text{sta}_1(F)$ ,  $A, F \models \text{f.frb}_1(p, q)$  iff  $l_1(p) > l_1(p, q)$ .

*Proof.* By definition,  $l_1(p) \geq l_1(p, q)$  for all  $p, q$ . We have  $l_1(p) > l_1(p, q)$  iff  $\exists p', q'$  such that  $A \models p' @ p \wedge A, F \models \text{f.acc}(p', q') \wedge F \not\models q' @ q$ . Since  $q \in \text{sta}_1(F)$  by assumption, this is equivalent to  $A, F \models \text{f.frb}_1(p, q)$ .  $\square$

**Theorem 23.** For stepwise tree automata with  $\epsilon$ -rules  $A$  and deterministic factorized tree automata  $F$  over the same signature, inclusion  $L(A) \subseteq L(F)$  can be decided in time  $O(|A| \cdot |F|)$ .

*Proof.* The algorithm first computes  $\text{lfp}(D_1(A, F))$  in time  $O(|A| \cdot |F|)$ . It returns false if the fixed point contains  $\text{fail}_0$ ,  $\text{fail}_1$ , or  $\text{fail}_2$ . Otherwise, it computes the values of all counters  $l_i(p)$  and  $l_i(p, q)$ . If  $l_i(p) > l_i(p, q)$  for some  $\text{acc}(p, q) \in \text{lfp}(D_1(A, F))$  then the algorithm returns false, otherwise true. All these steps can be performed in time  $O(|A| \cdot |F|)$  as argued above.  $\square$

$$\begin{array}{c}
\text{(f.frb}_1^c) \frac{A \models p_1 @ p_2 \quad q_2 \in \text{sta}_2(F)}{\text{f.frb}_1^c(p_1, Q_1^F(q_2)) :- \text{f.acc}(p_2, q_2)}. \\
\text{(f.frb}_2^c) \frac{A \models p_1 @ p_2 \quad q_1 \in \text{sta}_1(F)}{\text{f.frb}_2^c(p_2, Q_2^F(q_1)) :- \text{f.acc}(p_1, q_1)}.
\end{array}$$

Figure 15:  $D_2(A, F)$  extends  $D_1(A, F)$  with clauses for  $\text{f.frb}_i^c$ .

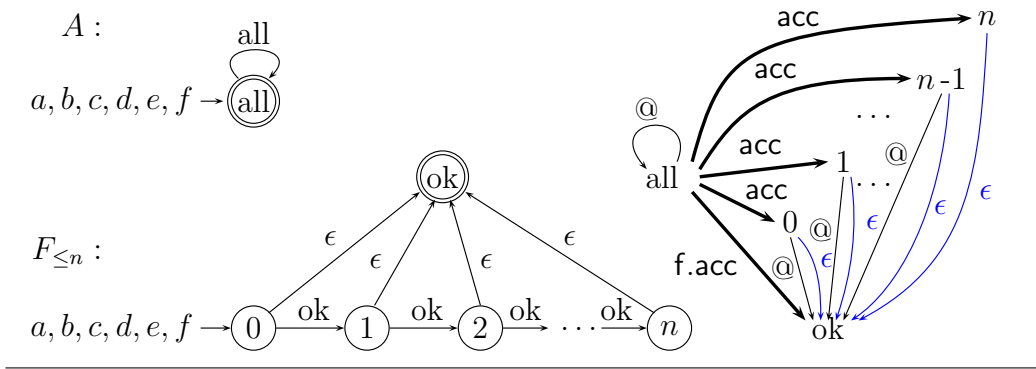


Figure 16: Example run of the algorithm:  $\text{f.frb}_1^c(\text{all}, \{0, \dots, n-1\})$  is inferred.

#### 4.4. Efficient Algorithm

We present a more efficient method to compute the values of the counters, that is similar to the non-factorized case. We use new predicates  $\text{f.frb}_i^c$  which account for complementation with respect to sort  $i$ , where  $i \in \{1, 2\}$ :

$$A, F \models \text{f.frb}_i^c(p, Q) \Leftrightarrow_{df} \forall q \in \text{sta}_i(F) \setminus Q, A, F \models \text{f.frb}_i(p, q)$$

Datalog program  $D_2(A, F)$  in Figure 15 infers  $\text{f.frb}_i^c(p, Q)$  literals. It extends  $D_1(A, F)$  by clauses from two further transformation rules (f.frb<sub>i</sub><sup>c</sup>).

**Lemma 24.** *If  $D = D_2(A, F)$  then  $l_i(p) = \sum_{Q \subseteq \text{sta}_i(F)} \text{lf}p^\#(D)(\text{f.frb}_i^c(p, Q))$  and  $l_i(p, q) = \sum_{Q \subseteq \text{sta}_i(F), q \in Q} \text{lf}p^\#(D)(\text{f.frb}_i^c(p, Q))$ .*

*Proof.* Similar to the proof of Lemma 12. □

**Lemma 25.**  $D_2(A, F)$  can be computed in time  $O(|A| \cdot |F|)$  from  $A$  and  $F$ .

*Proof.* The proof works as for  $D_2(A, B)$  in Lemma 10. The grouping clauses produced by (f.frb<sub>i</sub><sup>c</sup>) can be rewritten in analogy to those for (frb<sup>c</sup>) before. The sets  $Q_i^F$  are of size  $O(|F|)$  and occur at most  $O(|A|)$  times in rules (f.frb<sub>i</sub><sup>c</sup>), thus the overall size of the clauses produced by this rule is in  $O(|A| \cdot |F|)$ , too. The analysis for the remaining rules is straightforward. □

An example for the algorithm is given in Figure 16. Automaton  $A$  given there recognizes all trees, and the factorized tree automata  $F_{\leq n}$  all Curried encodings of unranked trees (see Section 5 for the definitions) with at most  $n$  children per node. Datalog program  $D_2(A, F_{\leq n})$  infers the literals  $\text{f.frb}_1^c(\text{all}, \{0, \dots, n-1\})$  and  $\text{f.frb}_2^c(\text{all}, \{\text{ok}\})$  as illustrated on the right of the figure. The first implies  $A, b(F_{\leq n}) \models \text{f.frb}_1(\text{all}, n)$  and thus  $A, b(F_{\leq n}) \models \text{fail}_1$ . The second is a tautology since  $\text{sta}_2(F_{\leq n}) = \{\text{ok}\}$ .

Compared to the non-factorized case, our algorithm cannot infer  $\text{frb}^c(p, Q)$  literals efficiently any more. This would require to apply  $\epsilon$ -rules over and over, spoiling our time complexity of  $O(|A| \cdot |F|)$ . Instead, our algorithm infers  $\text{f.frb}_i^c(p, Q)$  literals and combines them with  $\text{f.acc}(p, q)$  literals from  $(\text{f.frb}_i^c)$  in Figure 15. Epsilon-rules are used for inferring the  $\text{f.acc}(p, q)$  literals (see Figure 12). Besides, they only serve in transformation rules  $(\text{frb}_{/ia})$  from Figure 14, which deal with cases where evaluation stops in some state that cannot be converted to the required sort by any  $\epsilon$ -rule.

#### 4.5. Early Failure Detection

We show how to check for  $\text{fail}_1$  on the fly. We assume that all counters  $l_i(p)$  and  $l_i(p, q)$  are always up to date. As in the non-factorized case, we assume that literals with predicates  $\text{acc}$  are inferred with the lowest priority (also lower than  $\text{f.acc}$ ).

The *on-the-fly algorithm* works as follows. It computes  $\text{lf}_p(D_2(A, F))$  while returning false once  $\text{fail}_0$ ,  $\text{fail}_1$  or  $\text{fail}_2$  is produced. The counters  $l_i(p)$  and  $l_i(p, q)$  are always kept up-to-date, i.e., increased when some literal  $\text{f.frb}_i^c(p, Q)$  is inferred by some further clause. For all newly inferred literals  $\text{acc}(p, q)$ , it checks whether some literal  $\text{f.frb}_i(p, q)$  was produced before, where  $i \in \{1, 2\}$  and  $r \in \text{sta}_i(F) \setminus Q$  with  $q \xrightarrow{F}^{\leq 1} r$ . The existence of a literal  $\text{f.frb}_i(p, q)$  is reduced to checking whether  $l_1(p) > l_1(p, r)$  or  $l_2(p) > l_2(p, r)$ . If so,  $\text{fail}_1$  is raised and the algorithm returns false. Otherwise, it continues with the fixed point computation. Testing the counters on the fly is sufficient, as shown by the following lemma.

**Lemma 26.** *The on-the-fly algorithm detects  $\text{fail}_1$  if  $\text{acc}(p, q)$  is inferred before some literal  $\text{f.frb}_i^c(p, Q)$ , where  $r \notin Q$  is the unique state with  $q \xrightarrow{F}^{\leq 1} r$ .*

*Proof.* We consider the case  $i = 1$  only, which is sufficient by symmetry. Let  $\text{f.frb}_1^c(p_1, R_1)$  be inferred after  $\text{acc}(p_1, q_1)$ , where  $q_1 \xrightarrow{F}^{\leq 1} r_1$  and  $r_1 \in \text{sta}_1(F) \setminus R_1$ . See Figure 17 for illustration. Literal  $\text{f.frb}_1^c(p_1, R_1)$  is justified by some literal  $\text{f.acc}(p_2, r_2)$  added before, and the clause below where  $R_1 = Q_1^F(r_2)$ . Furthermore,  $\text{f.acc}(p_2, r_2)$  stems from some literal  $\text{acc}(p_2, q_2)$  and the second clause:

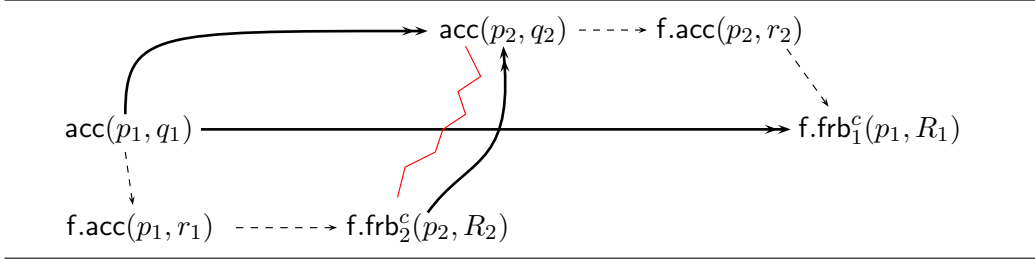


Figure 17: Early failure detection for  $i = 1$ .

$$\frac{A \models p_1 @ p_2 \quad r_2 \in sta_2(F)}{f.frb_1^c(p_1, R_1) :- f.acc(p_2, r_2)} \quad \frac{p_2 \in sta(A) \quad q_2 \xrightarrow{F}^{\leq 1} r_2}{f.acc(p_2, r_2) :- acc(p_2, q_2)}$$

Due to the lowest priority of `acc` literals again, `acc(p1, q1)` must be inferred before `acc(p2, q2)`. The following clauses can be applied where  $R_2 = Q_2^F(r_1)$ :

$$\frac{p_1 \in sta(A) \quad q_1 \xrightarrow{F}^{\leq 1} r_1}{f.acc(p_1, r_1) :- acc(p_1, q_1)} \quad \frac{A \models p_1 @ p_2 \quad r_1 \in sta_1(F)}{f.frb_2^c(p_2, R_2) :- f.acc(p_1, r_1)}$$

Since  $r_1 \in Q_1^F(r_2)$  if and only if  $r_2 \in Q_2^F(r_1)$ , it follows from  $r_1 \notin R_1$  that  $r_2 \notin R_2$ . Thus, `acc(p2, q2)` and `f.frb2c(p2, R2)` in  $lfp(D_2(A, F))$  raise inclusion failure `fail1`. By priority, `f.frb2c(p2, R2)` is added before `acc(p2, q2)`, so this failure is properly detected by the incremental algorithm.  $\square$

As in the non-factorized case, we can turn this algorithm incremental with respect to adding epsilon edges to  $A$ . We can thus test inclusion in deterministic factorized automata as efficiently as for the non-factorized case.

## 5. DTDs and Factorized Tree Automata for Unranked Trees

We lift our inclusion test to factorized tree automata interpreted over unranked trees, so that it becomes applicable to deterministic DTDs. Factorization is essential for efficiency here.

### 5.1. Factorized Tree Automata for Unranked Trees

An unranked signature  $\Sigma$  is a finite set of symbols (without arity restrictions). The set  $T_\Sigma^u$  of unranked trees over  $\Sigma$  is the least set containing all tuples  $a(t_1, \dots, t_n)$  where  $a \in \Sigma$ ,  $t_1, \dots, t_n \in T_\Sigma^u$  and  $n \geq 0$ .

Currying carries over literally from ranked to unranked trees. This yields the bijective function  $curry : T_\Sigma^u \rightarrow T_{\Sigma @}$ , which satisfies  $curry(a(t_1, \dots, t_n)) = a @ curry(t_1) @ \dots @ curry(t_n)$  for all unranked trees  $a(t_1, \dots, t_n) \in T_\Sigma^u$ . For instance,  $curry(a(b, c, d(e))) = a @ b @ c @ (d @ e)$ . Subtrees of  $a(b, c, d(e))$  are encoded as subtrees on the right of `@` such as  $d @ e$ .

Subtrees on the left of @ encode rooted hedges such as  $a@b@c$  that are subject to extension to the right. This semantic difference motivated different sorts for hedges and unranked trees already in the automata notions of Raeymaekers [14] or Neumann & Seidl [4].

We can use factorized tree automata  $A$  over stepwise signatures  $\Sigma$  to recognize languages of unranked trees  $L^u(A) = \{t \in T^u(\Sigma) \mid \text{curry}(t) \in L(A)\}$ . We obtain the following corollary from Theorem 23.

**Corollary 27.** *Let  $A$  be a stepwise tree automaton and  $F$  a deterministic factorized tree automaton over the same signature  $\Sigma$ . Language inclusion  $L^u(A) \subseteq L^u(F)$  can be decided in time  $O(|A| \cdot |F|)$  independently of  $|\Sigma|$ .*

The tree automaton  $A$  can also be chosen to be a hedge automaton, whose horizontal languages are defined by non-deterministic finite word automata (nFAs). Hedge automata  $H$  over  $\Sigma$  have rules of the form  $a(C) \rightarrow q$  where  $a \in \Sigma$ ,  $q \in \text{sta}(H)$ , and  $C$  is an nFA with signature  $\text{sta}(H)$ . Such hedge automata are called NFHAs by Comon *et al.* [1] and UTAs by Martens & Niehren [13]. They can be translated in linear time to stepwise tree automata with  $\epsilon$ -rules [13]. A hedge automaton is called deterministic if all its nFAs are deterministic (dFAs) and  $L(q_1) \cap L(q_2) = \emptyset$  for all two rules  $a(C_1) \rightarrow q_1$  and  $a(C_2) \rightarrow q_2$  in  $\text{rul}(H)$ . This is only a pseudo-notion of determinism. It is mapped to unambiguity of stepwise tree automata. As a consequence, we cannot choose  $F$  to be a deterministic hedge automaton.

## 5.2. Deterministic DTDs

We convert deterministic DTDs  $D$  to deterministic factorized tree automata for unranked trees in time  $O(|\Sigma| \cdot |D|)$ , so that we can reuse our algorithm for testing inclusion of stepwise tree automata in deterministic DTDs. Here, factorization avoids the quadratic blowup. When translating into stepwise tree automata [13], the number of rules may become quadratic, while the number of states is preserved. The problem is the implicit elimination of  $\epsilon$ -rules.

A DTD  $D$  with elements in a set  $\Sigma$  is a function mapping letters  $a \in \Sigma$  to regular expressions  $e$  over  $\Sigma$ , in which case we write  $a \rightarrow_D e$ . One of these elements is the distinguished start symbol. Let  $L(e) \subseteq \Sigma^*$  be the word language defined by  $e$ . The language  $L_a(D) \subseteq T_\Sigma^u$  of elements  $a$  of a DTD  $D$  is the smallest set of unranked trees such that:

$$L_a(D) = \{a(t_1, \dots, t_n) \mid a \rightarrow_D e, a_1 \dots a_n \in L(e), t_i \in L_{a_i}(D) \text{ for } 1 \leq i \leq n\}$$

The language of a DTD  $D$  is  $L(D) = L_a(D)$  where  $a$  is the start symbol of  $D$ . The size of  $D$  is the total number of symbols in the regular expressions of  $D$ .



doc	$\rightarrow$ block <sup>+</sup>	<!ELEMENT doc (block+)>
block	$\rightarrow$ text (link text)? + link text?	<!ELEMENT block (text,(link ,text)?  link ,text?)>
text	$\rightarrow$ $\epsilon$	<!ELEMENT text (#PCDATA)>
link	$\rightarrow$ $\epsilon$	<!ELEMENT link (#PCDATA)>

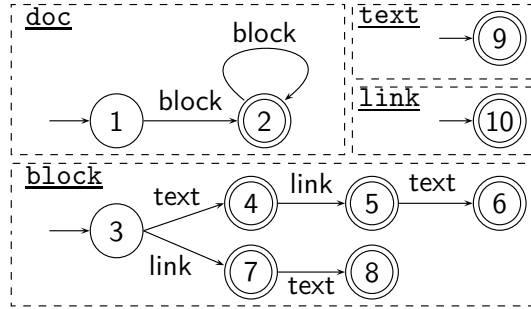


Figure 18: An example DTD and the corresponding Glushkov automata.

An example with its corresponding XML syntax is given in Figure 18. The set of elements of  $D$  is  $\Sigma = \{\text{doc}, \text{block}, \text{text}, \text{link}\}$ , of which the element  $\text{doc}$  is the start symbol. The regular expression for  $\#PCDATA$  recognizes only the empty word.

A DTD is deterministic if all its regular expressions are one-unambiguous, as required by the W3C. This is equivalent to say that all corresponding Glushkov automata are deterministic [12]. Glushkov automata are nFAs from regular expressions as usual except that  $\epsilon$  rules are eliminated on the fly whenever they appear. The precise definition is outside the scope of this article, but an example is given in Figure 18.

**Theorem 28 (Brüggemann-Klein [28]).** *The collection of Glushkov automata for a deterministic DTD  $D$  over  $\Sigma$  can be computed in time  $O(|\Sigma| \cdot |D|)$ .*

Note that the construction of the Glushkov automaton of a regular expression  $e$  over alphabet  $\Sigma$  may take time  $O(|\Sigma| \cdot |e|^2)$  in the general case. Intuitively, the square factor is raised by eliminating occurring  $\epsilon$ -rules on the fly. In the case of a one-unambiguous regular expression, the resulting Glushkov automaton is deterministic. The construction time is bounded by its size and thus in  $O(|\Sigma| \cdot |e|)$  due to determinism.

We transform the collection of Glushkov automata for a deterministic DTD  $D$  into a single factorized tree automaton  $F$  as follows. The set of states of sort 1 of  $F$  is the disjoint union of the states of the Glushkov

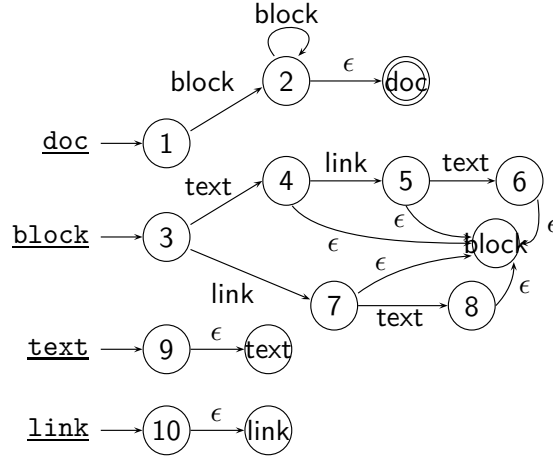


Figure 19: A representation of the deterministic factorized tree automaton for the DTD in Figure 18. Alphabet  $\Sigma$  is  $\{\underline{\text{doc}}, \underline{\text{block}}, \underline{\text{text}}, \underline{\text{link}}\}$ , states of sort 1 are  $\{1, \dots, 10\}$  and states of sort 2 are  $\{\underline{\text{doc}}, \underline{\text{block}}, \underline{\text{text}}, \underline{\text{link}}\}$ . A constant rule is for instance  $\underline{\text{doc}} \rightarrow 1$ , a binary rule  $2@block \rightarrow 2$  and an  $\epsilon$ -rule  $2 \xrightarrow{\epsilon} \underline{\text{doc}}$ .

automata. The states of sort 2 of  $F$  are the elements of  $D$ . For every element  $a$ , we connect all final states  $q$  of its Glushkov automaton to the state  $a$ , i.e.,  $q \xrightarrow{\epsilon} a \in \text{rul}(F)$ . The only final state of  $F$  is the start symbol of the DTD  $D$ . The result is an nFA that represents a factorized tree automaton, as for instance in Figure 19. This needs time of at most  $O(|\Sigma| \cdot |D|)$ . For every  $a \in \Sigma$ , there is a rule  $a \rightarrow q \in \text{rul}(F)$  for the unique initial state  $q$  of the Glushkov automaton of  $a$ . For every transition  $q \xrightarrow{a} q'$  of one of the Glushkov automata, we add a rule  $q@a \rightarrow q' \in \text{rul}(F)$ .

Note that  $F$  is deterministic as a factorized automaton. The  $\epsilon$ -free part of  $F$  is deterministic since all Glushkov automata are, thus establishing  $\mathbf{d}_0$ . Let  $q$  be a state of the Glushkov automaton for some letter  $a$ . The only state of sort 1 that  $q$  can reach by  $\epsilon$ -edges in  $F$  is  $a$  and the only state of sort 2 is  $q$  itself. All other states of  $F$  are elements of  $\Sigma$ , which have no outgoing  $\epsilon$ -edges, thus establishing  $\mathbf{d}_1$ . Note that the size of the example automaton would grow quadratically, when eliminating  $\epsilon$ -edges.

**Theorem 29.** *Deterministic DTDs  $D$  over  $\Sigma$  can be translated in time  $O(|\Sigma| \cdot |D|)$  to bottom-up deterministic factorized tree automata recognizing the same language.*

*Proof.* The translation of a collection of Glushkov automata of a DTD to a factorized automaton is in linear time. It is easy to check that it preserves the languages of unranked trees. The theorem thus follows from Theorem 28 by Brüggemann-Klein.  $\square$

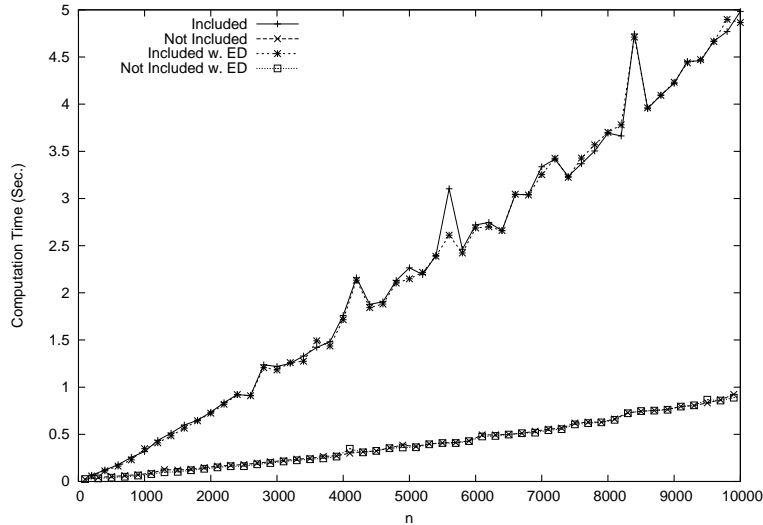


Figure 20: Computation time for testing  $L(Mult_n) \subseteq L(Mult_{200})$ .

**Corollary 30.** *Language inclusion of hedge automata  $A$  over  $\Sigma$  with horizontal languages defined by finite word automata in deterministic DTDs  $D$  with elements in  $\Sigma$  can be decided in time  $O(|A| \cdot |\Sigma| \cdot |D|)$ .*

*Proof.* From Corollary 27 and Theorem 29. □

## 6. Experiments

We have implemented the inclusion algorithm in Objective CAML, and have integrated it into a system for schema-guided learning of queries in XML trees [9]. In the first set of experiments, we consider inclusion tests for synthetic automata. Then we consider inclusion tests between automata and DTDs coming from realistic tasks in query learning.

**Experiment 1.** We modify the sizes of automata  $A$  and  $F$  when testing inclusion of  $L(A)$  in  $L(F)$ . For this, we define  $Mult_n$  as the minimal deterministic automaton for the language of trees of the form  $f(a, \dots, a)$  where the number of  $a$ -leaves is a multiple of  $n$ . The first problem is to test inclusion of  $L(Mult_n)$ ,  $n$  varying from 100 to 10000 with a 100-increment, into the minimal deterministic factorized automaton recognizing  $L(Mult_{200})$ . It should be noted that inclusion holds when  $n/100$  is even. The second problem is to test inclusion of  $L(Mult_{400})$  into  $L(Mult_n)$  with  $n$  varying from 10 to 500 with a 10-increment. It should be noted that inclusion holds when  $400/n$  is an integer.

We estimate the computation time for inclusion tests with and without early detection of inclusion failure (ED). We distinguish whether inclusion

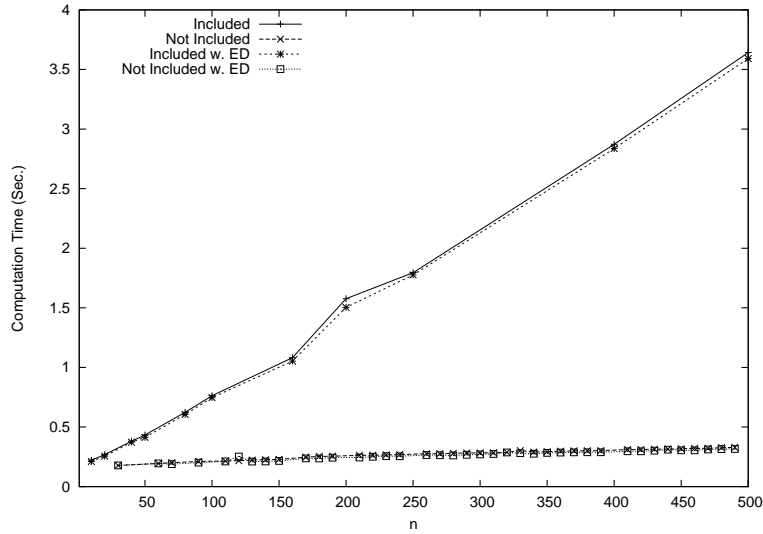


Figure 21: Computation time for testing  $L(Mult_{400}) \subseteq L(Mult_n)$ .

holds or not. Results are shown in Figure 20 for the first problem and in Figure 21 for the second problem.

It can be verified that the computation time of testing  $L(A) \subseteq L(F)$  is linear in the size of the automaton  $A$  and in the size of automaton  $F$ . This confirms the theoretical results on complexity. It can also be seen that the computation time is greater when inclusion holds. Otherwise, the computation time is lower since concurrent failure detection applies. In this experiment, there are no failures of type `fail1`, so we do not use early failure detection. The gain is obtained by checking `fail2` concurrently, so that product automaton does not need to be computed entirely.

**Experiment 2.** In order to verify the usefulness of early detection of failures of type `fail1` (ED), we consider another example. We define  $Mult2_n$  to be the minimal deterministic automaton for the language of trees of the form  $g(f(a, \dots, a))$ , where the number of  $a$ -leaves is a multiple of  $n$ . The problem is to test the inclusion of  $L(Mult2_n)$ , where  $n$  varies from 100 to 10000 with a 100-increment, into the minimal deterministic factorized automaton recognizing  $L(Mult2_{200})$ . The computation times are shown in Figure 22.

It can be noted that when inclusion does not hold, computation time is five times faster than for other cases. This is because, for these inclusion tests, inclusion failure comes from `fail1`. Thus early failure detection allows to decrease dramatically the computation time. It can also be noted that the computation time is similar for cases where inclusion is verified (with or without early detection), and non-inclusion cases without early detection.

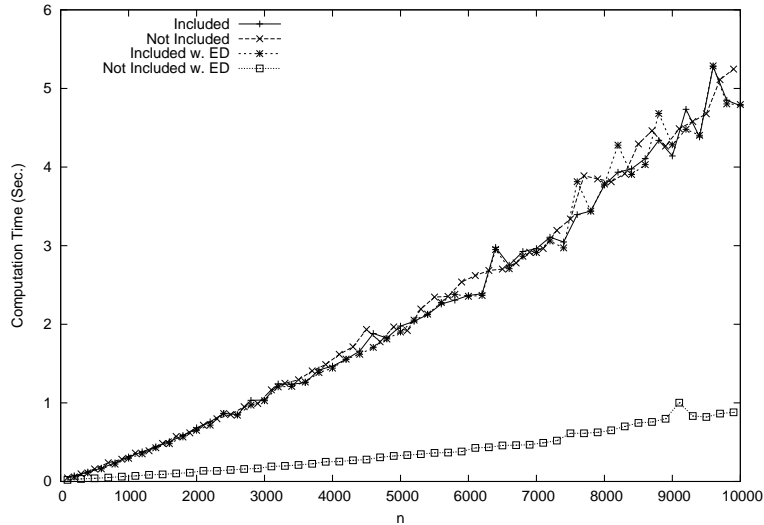


Figure 22: Average computation time for testing  $L(Mult2_n) \subseteq L(Mult2_{200})$ .

This is because, without early failure detection, computing a failure  $fail_1$  implies the need to compute the whole product automaton.

**Experiment 3.** We now consider real-world data sets from the query induction problem. In the learning algorithm defined by Champavère *et al.* [9], an initial automaton is computed and is iteratively refined by merging states. A merge is accepted only if the language recognized by the new automaton still satisfies a given schema or DTD. Consequently, inclusion tests are done frequently. We compare the overall computation time of learning sessions where inclusion tests are done with or without early failure detection. We use the the transitional DTD of XHTML and the query learning benchmarks Okra, Bigbook, Google and Yahoo, each of them with an increasing size of inputs. Results are shown in Figure 23.

It appears that the learning algorithm operates about twice as fast with early detection of failure  $fail_1$  than without in all benchmarks. This indicates that  $fail_1$  occurs frequently in practice and that early detection improves efficiency a lot. More generally, it shows that the inclusion algorithm presented here can be used for real-world problems. In Champavère *et al.* [9], we have shown that introducing inclusion tests does not increase computation time while avoiding useless state merges, thus improving the query induction algorithm.

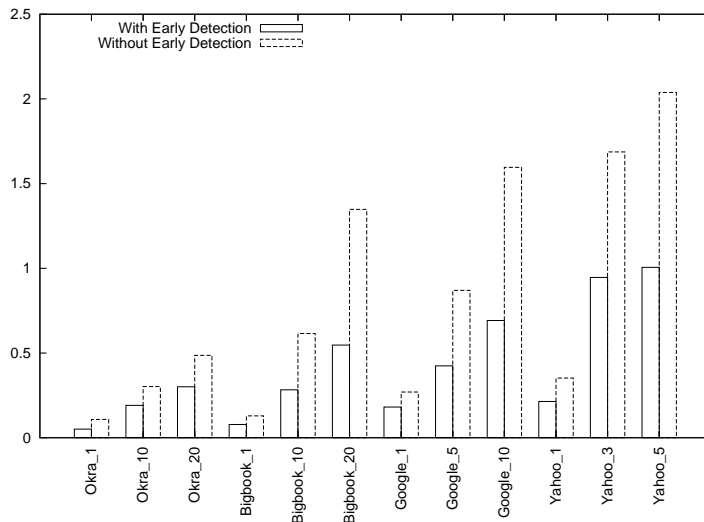


Figure 23: Average computation time (in Sec.) of learning sessions with and without early detection of failure `fail1`.

## 7. Top-Down Determinism

We show how to test inclusion for top-down deterministic tree automata by reduction to inclusion testing in deterministic word automata. This reduction should be folklore (but we are not aware of any reference). For our precise complexity analysis, however, we use Theorem 9 for the case of words, seen as trees over monadic signatures.

Thereby we obtain an efficient test for inclusion in deterministic EDTDs with restrained competition (and thus for schema definitions in XML Schema), as these can be translated to top-down deterministic tree automata with respect to Rabin’s firstchild-nextsibling encoding of unranked trees (see below). More precisely, we can test language inclusion for tree automata  $A$  recognizing firstchild-nextsibling encodings of unranked trees in deterministic restrained competition EDTDs  $D$  in time  $O(|A| \cdot |\Sigma| \cdot |D|)$ . We do not know whether the analogous result holds for automata  $A$  recognizing curried encodings.

A restriction of essentially the same algorithm for testing inclusion between two deterministic restrained competition EDTDs was presented earlier by Martens, Neven, Schwentick, & Bex [18] (see Section 10 of the reference). Their presentation, however, does not rely on top-down deterministic tree automata as an intermediate step, and no precise complexity analysis is given.

### 7.1. Top-Down Deterministic Tree Automata for Ranked Trees

A tree automaton  $A$  over a ranked signature  $\Sigma$  is *top-down deterministic* if for all symbols  $f \in \Sigma$  of arity  $n$  and states  $p \in \Sigma$  there are no two different rules  $f(p_1, \dots, p_n) \rightarrow p$  and  $f(p'_1, \dots, p'_n) \rightarrow p$  in  $\text{rul}(A)$ .

**Proposition 31.** *Let  $\Sigma$  be a ranked signature, and  $A$  and  $B$  be tree automata over  $\Sigma$ . If  $B$  is top-down deterministic, then we can decide language inclusion  $L(A) \subseteq L(B)$  in time  $O(|A| \cdot |B|)$ .*

We base the algorithm on the well-known fact that tree languages recognized by top-down deterministic tree automata are path-closed [30, 1]. The standard example for a non-path-closed regular language is  $L_0 = \{f(a, a), f(b, b)\}$  where  $a \neq b$ . For the sake of completeness, let us recall the definitions. The set of paths of a tree  $t \in T_\Sigma$  is the subset of words  $\text{paths}(t) \subseteq (\Sigma \cup \mathbb{N})^*$  defined as follows:

$$\text{paths}(a) = a, \quad \text{and} \quad \text{paths}(f(t_1, \dots, t_n)) = \{fiw \mid 1 \leq i \leq n, w \in \text{paths}(t_i)\}.$$

For instance  $\text{paths}(L_0) = \{f1a, f2a, f1b, f2b\}$ . The path closure of a tree language  $L \subseteq T_\Sigma$  is the set of all trees that contain only paths of trees in  $L$ :

$$\text{path-clos}(L) = \{t \mid \text{paths}(t) \subseteq \text{paths}(L)\}$$

We call  $L$  path-closed if  $L = \text{path-clos}(L)$ . For instance  $\text{path-clos}(L_0) = L_0 \cup \{f(a, b), f(b, a)\}$ , so  $L_0$  is indeed not path-closed.

**Lemma 32.** *If  $L_2$  is path-closed then  $L_1 \subseteq L_2$  iff  $\text{paths}(L_1) \subseteq \text{paths}(L_2)$ .*

*Proof.* Note that we do not assume  $L_1$  to be path-closed. The implication from left to right is trivial. For the inverse, assume  $\text{paths}(L_1) \subseteq \text{paths}(L_2)$ . If  $t_1 \in L_1$  then  $\text{paths}(t_1) \subseteq \text{paths}(L_1) \subseteq \text{paths}(L_2)$ . Thus  $t_1 \in \text{path-clos}(L_2)$ , and this set is equal to  $L_2$  by assumption of path-closeness.  $\square$

*Proof of Proposition 31.* For every tree automaton  $A$  over  $\Sigma$ , we construct an nFA  $P(A)$  over a finite subset of  $\Sigma \uplus \mathbb{N}$  such that  $L(P(A)) = \text{paths}(L(A))$ . The rules of  $P(A)$  are defined as follows:

$$\begin{aligned} \text{rul}(P(A)) &= \{p \xrightarrow{f_i} p_i \mid f(p_1, \dots, p_n) \rightarrow p \in \text{rul}(A), 1 \leq i \leq n\} \\ &\cup \{a \rightarrow p \mid a \rightarrow p \in \text{rul}(A)\} \\ &\cup \{p \xrightarrow{\epsilon} p' \mid p \xrightarrow{\epsilon} p' \in \text{rul}(A)\} \end{aligned}$$

The first kind of rules reads two letters at the same time, but can be easily rewritten into two rules reading each a single letter. Clearly, the construction of  $P(A)$  is in time  $O(|A|)$ . Furthermore,  $P(A)$  is deterministic iff  $A$  is top-down deterministic.

Given two tree automata  $A, B$  over  $\Sigma$  such that  $B$  is top-down deterministic, we can decide language inclusion between  $A$  and  $B$  by testing language inclusion for  $P(A)$  and  $P(B)$ . Since  $P(B)$  is deterministic this can be done in time  $O(|P(A)| \cdot |P(B)|)$  independently of the alphabet by Theorem 9, which is in time  $O(|A| \cdot |B|)$  independently of the signature.  $\square$

### 7.2. Restrained Competition EDTDs

We test inclusion in restrained competition EDTDs, for tree automata recognizing unranked trees modulo the firstchild-nextsibling encoding of unranked trees. It is obtained by encoding restrained competition EDTDs to top-down deterministic tree automata with respect to this binary encoding.

An extended DTD (EDTD)  $D$  over a signature  $\Sigma$  consists of a finite set of states  $sta(D) \subseteq \Sigma \times \mathbb{N}$ , a subset of start states  $start(D) \subseteq sta(D)$ , and a collection of rules given by a function mapping states  $q \in sta(D)$  to regular expressions  $e$  over  $sta(D)$ , in which case we write  $q \rightarrow_D e$ . The language  $L_q(D) \subseteq T_\Sigma^u$  of a state  $q \in sta(D)$  is the smallest set of unranked trees such that if  $q = (a, i)$  for some  $a \in \Sigma, i \in \mathbb{N}$  then:

$$L_q(D) = \{a(t_1, \dots, t_n) \mid q \rightarrow_D e, q_1 \dots q_n \in L(e), t_i \in L_{q_i}(D) \text{ for } 1 \leq i \leq n\}$$

The language of an EDTD is  $L(D) = \cup_{q \in start(D)} L_q(D)$ . The size of  $D$  is the total number of symbols in the regular expressions of  $D$ . Essentially, EDTDs are the same as hedge automata with regular expressions for defining horizontal languages. They can thus recognize all regular languages of unranked trees.

An EDTD  $D$  is *restrained competition* if it has a unique start state and for all regular expressions  $q \rightarrow_D e$  of  $D$  there exist no two different states  $(a, n_1), (a, n_2) \in sta(D)$  and words  $u, v_1, v_2 \in sta(D)^*$  such that  $u(a, n_1)v_1, u(a, n_2)v_2 \in L(e)$ . Also, as for DTDs, we call a restrained competition EDTD *deterministic* if all its regular expressions are one-unambiguous and if it has at most one start state.

Restrained competition EDTDs are strictly more expressive than deterministic DTDs. On the other hand, they are more restrictive than the class of regular languages, in order to permit the typing of all nodes of an XML document in 1-pass streaming manner [31, 18]. Consider for instance the regular language of unranked trees  $L_1 = \{a(a)\}$ . It cannot be recognized by any DTD, since it contains a tree with two types of  $a$ -nodes that need to be



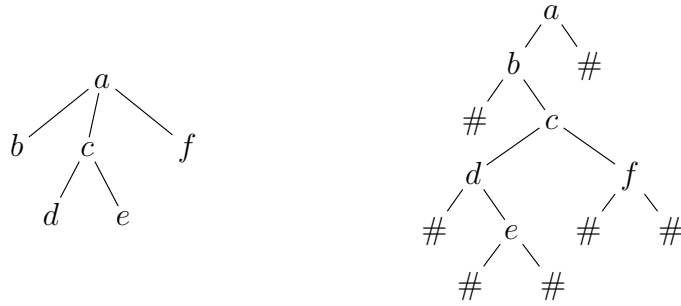


Figure 24: The tree  $a(b, c(d, e), f)$  and its firstchild-nextsibling encoding  $a(b(\#, c(d(\#, e(\#, \#)), f(\#, \#))), \#)$ .

distinguished. Language  $L_1$  can be recognized by a restrained competition EDTD with two states  $(a, 1)$  and  $(a, 2)$  for the two types of  $a$ -nodes. The start state is  $(a, 1)$  and the rules are as follows:

$$(a, 1) \rightarrow (a, 2), \quad (a, 2) \rightarrow \epsilon$$

This example illustrates that we cannot translate restrained competition EDTDs to bottom-up deterministic stepwise tree automata in linear time (in contrast to the case of deterministic DTDs). The naive approach would be to permit tree automaton rules  $a \rightarrow (a, 1)$  and  $a \rightarrow (a, 2)$  but these violate bottom-up determinism. The problem is that the type of an  $a$ -node is only determined once knowing the type of its parent, so we have to try out all choices in a bottom-up deterministic manner.

Let  $\Sigma^\# = \Sigma \uplus \{\#\}$  be the ranked signature with a single constant  $\#$  and a collection of binary function symbols  $a \in \Sigma$ . Rabin's firstchild-nextsibling encoding  $fcns$  of an unranked tree  $t \in T^u(\Sigma)$  is a binary tree in  $T_{\Sigma^\#}$  (see, e.g., Gottlob & Koch [29]), for instance,  $fcns(a(b, c(d, e), f)) = a(b(\#, c(d(\#, e(\#, \#)), f(\#, \#))), \#)$  as illustrated in Figure 24. A tree automaton  $A$  over  $\Sigma^\#$  recognizes unranked trees modulo this other binary encoding, so its unranked tree language is  $L^u(A) = \{t \in T^u(\Sigma) \mid fcns(t) \in L(A)\}$ .

**Lemma 33.** *For all deterministic restrained competition EDTDs  $D$  over  $\Sigma$ , we can compute a top-down deterministic tree automaton  $B$  over  $\Sigma^\#$  with the same unranked tree language  $L^u(B) = L(D)$  in time  $O(|\Sigma| \cdot |D|)$ .*

*Proof.* We first compute the collection of Glushkov automata  $G_q$  for all regular expressions  $e$  such that  $q \rightarrow_D e$ . Since  $D$  is deterministic, all regular expressions  $e$  are unambiguous, so that all Glushkov automata  $G_q$  are dFAs of overall size  $O(|\Sigma| \cdot |D|)$  by Theorem 28. The alphabets of  $G_q$ s is

$sta(D) \subseteq \Sigma \times \mathbb{N}$ . Without restriction of generality, we can assume that  $G_q$  is productive. For productive  $G_q$ , restrained competition of  $D$  implies that there are no two rules  $p \xrightarrow{(a,i)} p' \in rul(G_q)$  and  $p \xrightarrow{(a,j)} p'' \in rul(G_q)$  with  $i \neq j$ . Determinism of  $G_q$  implies that:

(\*) for all  $p \in sta(G_q)$  and letters  $a \in \Sigma$  there is at most one pair  $(i, p')$  such that  $p \xrightarrow{(a,i)} p' \in rul(G_q)$ .

From the collection  $(G_q)_{q \in sta(D)}$ , we build a tree automaton  $B$  over the signature  $\Sigma^\#$  with  $L^u(B) = L(D)$ . The states of automaton  $B$  are the elements in  $\{S, H\} \uplus (\uplus_{q \in sta(D)} sta(G_q))$ , of which only  $S$  is final, i.e., the state of the root. State  $H$  is the state of the hash symbol  $\#$  at the second child of the root. The rules in  $rul(B)$  are defined as follows, where  $I$  is the unique initial state of  $init(G_{(a^s, i^s)})$  and  $(a^s, i^s)$  the unique start state of  $D$ .

$$\frac{p \xrightarrow{(a,i)} p' \in rul(G_q) \quad init(G_{(a,i)}) = \{p''\}}{a(p'', p') \rightarrow p} \quad \frac{p \in fn(G_q)}{\# \rightarrow p} \quad \frac{true}{\begin{array}{l} a^s(I, H) \rightarrow S \\ \# \rightarrow H \end{array}}$$

Automaton  $B$  is top-down deterministic by (\*) and recognizes  $L(D)$ . The construction is in linear time in the size of the collection of Glushkov automata, so the overall construction requires time  $O(|\Sigma| \cdot |D|)$ .  $\square$

**Corollary 34.** *For tree automata  $A$  over  $\Sigma^\#$  and deterministic restrained competition EDTDs  $D$  over  $\Sigma$ , language inclusion  $L^u(A) \subseteq L(D)$  can be tested in time  $O(|A| \cdot |\Sigma| \cdot |D|)$ .*

*Proof.* We transform  $D$  into a top-down deterministic tree automaton  $B$  over  $\Sigma^\#$  that recognizes the same unranked tree language by Lemma 33. This takes time  $O(|\Sigma| \cdot |D|)$ , so the size of  $B$  is in  $O(|\Sigma| \cdot |D|)$  as well. We then test  $L(A) \subseteq L(B)$ , which can be done in time  $O(|A| \cdot |\Sigma| \cdot |D|)$  by Proposition 31, since  $B$  is top-down deterministic.  $\square$

## 8. Conclusion

We have presented new efficient algorithms for testing language inclusion in deterministic tree automata and XML schema definitions.

Our main contribution is an efficient inclusion test of tree automata  $A$  in bottom-up deterministic factorized tree automata  $B$  in time  $O(|A| \cdot |B|)$ . We have introduced early failure detection, which gives us the ability to turn our algorithm incremental with respect to adding epsilon rules to  $A$ . We

have implemented our inclusion test, given experimental evidence for its efficiency, and applied it in schema-guided query induction [9]. Incrementality considerably improves efficiency according to our experiments.

We have translated deterministic DTDs  $D$  to bottom-up deterministic factorized tree automata of size  $O(|\Sigma| \cdot |D|)$ . Our new notion of factorization is essential for efficiency here. As a corollary, we can check inclusion of stepwise tree automata  $A$  in deterministic DTDs  $D$  in time  $O(|A| \cdot |\Sigma| \cdot |D|)$ . Automata  $A$  can also be chosen to be hedge automata with nFAs for defining horizontal languages [1]. Such hedge automata are equivalent to EDTDs, since regular expressions can be translated to nFAs with  $\epsilon$ -rules in linear time. They can also be obtained from schema definitions in Relax NG.

We have presented a simpler inclusion test for inclusion in top-down deterministic tree automata in the case of ranked trees. This case is of interest for unranked trees, since deterministic DTDs and restrained competition EDTDs can be translated to top-down deterministic tree automata with respect to the firstchild-nextsibling encoding of unranked trees.

*Future Work.* One question on deterministic inclusion is left open, which is whether one can test inclusion of stepwise tree automata  $A$  in restrained competition EDTDs  $D$  over the same signature  $\Sigma$  in time  $O(|A| \cdot |\Sigma| \cdot |D|)$ . There are two difficulties here. First, it does not help to convert a stepwise tree automaton into a hedge automaton or a tree automaton operating on unranked trees modulo the firstchild-nextsibling encoding, since the known transformations introduce quadratic blowups. Second, we cannot represent the collection of Glushkov automata of a deterministic restrained competition EDTD by a deterministic stepwise tree automaton of the same size. The difference is that stepwise tree automata operate bottom-up while restrained competition EDTDs work top-down. We hope to solve this problem in future work by studying inclusion in deterministic streaming tree automata [32, 33, 4], which can operate in a mixed top-down and bottom-up manner.

#### *Acknowledgements*

We thank Sławek Staworko for having pointed out the alternative algorithm for inclusion checking for top-down deterministic tree automata and its relevance for DTDs.

We thank the anonymous reviewers for their valuable suggestions, which helped us a lot in order to simplify the presentation of the tedious counting arguments in the core of our algorithm and proofs.

The work was supported by the ANR project Codex.

## References

- [1] H. Comon, M. Dauchet, R. Gilleron, C. Löding, F. Jacquemard, D. Lugiez, S. Tison, M. Tommasi. Tree automata techniques and applications. Available online: <http://www.grappa.univ-lille3.fr/tata>, 1997, extended and revised in October 2007.
- [2] H. Seidl. Deciding equivalence of finite tree automata. *SIAM Journal on Computing* 19 (3) (1990) 424–437.
- [3] H. Seidl. Haskell overloading is DEXPTIME-complete. *Information Processing Letters* 52 (2) (1994) 57–60.
- [4] A. Neumann, H. Seidl. Locating matches of tree patterns in forests. *International Conference on Foundations of Software Technology and Theoretical Computer Science (FCTTCS)*, 1998, pp. 134–145.
- [5] S. D. Zilio, D. Lugiez. XML Schema, tree logic and sheaves automata. *International Conference on Rewriting Techniques and Applications (RTA)*, 2003, pp. 246–263.
- [6] S. Maneth, A. Berlea, T. Perst, H. Seidl. XML type checking with macro tree transducers. *Symposium on Principles of Database Systems (PODS)*, 2005, pp. 283–294.
- [7] T. Schwentick. Automata for XML—a survey. *Journal of Computer and System Science* 73 (3) (2007) 289–315.
- [8] T. Milo, D. Suciu, V. Vianu. Type checking XML transformers. *Journal of Computer and System Science* 1 (66) (2003) 66–97.
- [9] J. Champavère, R. Gilleron, A. Lemay, J. Niehren. Schema-guided induction of monadic queries. *International Colloquium on Grammatical Inference (ICGI)*, 2008, pp. 15–28.
- [10] J. Carme, J. Niehren, M. Tommasi. Querying unranked trees with step-wise tree automata. *International Conference on Rewriting Techniques and Applications (RTA)*, 2004, pp. 105–118.
- [11] G. J. Bex, F. Neven, J. V. den Bussche, DTDs versus XML Schema: a practical study. *International Workshop on the Web and Databases (WebDB)*, 2004, pp. 79–84.
- [12] A. Brüggemann-Klein, D. Wood. One-unambiguous regular languages. *Information and Computation* 142 (2) (1998) 182–206.

- [13] W. Martens, J. Niehren. On the minimization of XML schemas and tree automata for unranked trees. *Journal of Computer and System Science* 73 (4) (2007) 550–583.
- [14] S. Raeymaekers. Information extraction from web pages based on tree automata induction. Ph.D. thesis, Katholieke Universiteit Leuven (Jan. 2008).
- [15] J. W. Thatcher. Characterizing derivation trees of context-free grammars through a generalization of automata theory. *Journal of Computer and System Science* 1 (1967) 317–322.
- [16] A. Brüggemann-Klein, D. Wood, M. Murata. Regular tree and regular hedge languages over unranked alphabets: Version 1 (Apr. 07, 2001), <http://www.cs.ust.hk/tcsc/RR/2001-05.ps.gz>.
- [17] W. Martens, F. Neven, T. Schwentick. Which XML schemas admit 1-pass preorder typing? *International Conference on Database Theory*, 2005, pp. 68–82.
- [18] W. Martens, F. Neven, T. Schwentick, G. J. Bex. Expressiveness and complexity of XML Schema. *ACM Transactions on Database Systems* 31 (3) (2006) 770–813.
- [19] E. van der Vlist. *XML Schema*. O’Reilly Media, Inc., 2002.
- [20] E. van der Vlist. *Relax NG*. O’Reilly Media, Inc., 2003.
- [21] J. Champavère, R. Gilleron, A. Lemay, J. Niehren. Efficient inclusion checking for deterministic tree automata and DTDs. *International Conference on Language and Automata Theory and Applications (LATA)*, 2008, pp. 184–195.
- [22] W. Martens, F. Neven, T. Schwentick. Complexity of decision problems for simple regular expressions. *International Symposium on Mathematical Foundations of Computer Science (MFCS)*, 2004, pp. 889–900.
- [23] A. Tozawa, M. Hagiya. XML Schema containment checking based on semi-implicit techniques. *International Conference on Implementation and Application of Automata (CIAA)*, 2003, pp. 51–61.
- [24] S. Ceri, G. Gottlob, L. Tanca. What you always wanted to know about Datalog (and never dared to ask). *IEEE Transactions on Knowledge Data Engineering* 1 (1) (1989) 146–166.

- [25] E. Dantsin, T. Eiter, G. Gottlob, A. Voronkov. Complexity and expressive power of logic programming. *ACM computing surveys* 33 (3) (2001) 374–425.
- [26] G. Gottlob, E. Grädel, H. Veith. Datalog LITE: a deductive query language with linear time model checking. *ACM Transactions on Computational Logics* 3 (1) (2002) 42–79.
- [27] M. Minoux. LTUR: a simplified linear-time unit resolution algorithm for Horn formulae and computer implementation. *Information Processing Letters* 29 (1) (1988) 1–12.
- [28] A. Brüggemann-Klein. Regular expressions to finite automata. *Theoretical Computer Science* 120 (2) (1993) 197–213.
- [29] G. Gottlob, C. Koch. Monadic queries over tree-structured data. *IEEE Symposium on Logic in Computer Science (LICS)*, 2002, pp. 189–202.
- [30] J. Viràgh. Deterministic ascending tree automata I. *Acta Cybernetica* 5 (1) (1980) 33–42.
- [31] M. Murata, D. Lee, M. Mani. Taxonomy of XML schema languages using formal language theory. *ACM Transactions on Internet Technology* 5 (4) (2005) 660–704.
- [32] O. Gauwin, J. Niehren, Y. Roos. Streaming tree automata. *Information Processing Letters* 109 (1) (2008) 13–17.
- [33] R. Alur. Marrying words and trees. *Symposium on Principles of Database Systems (PODS)*, 2007, pp. 233–242.

Figure 25: Algorithm in pseudo-language for testing inclusion.

---

```

// Inputs:
// - A: productive stepwise tree automaton
// - F: deterministic factorized tree automaton over the same signature
// Output: true iff  $L(A) \subseteq L(F)$ 
fun inclusion(A,F)
  exception fail0 fail1 fail2
  << create counters >>
  << create literal collection >>
  << create agenda with priorities >>
  try
    // compute  $\text{lfp}(D_2(A,F))$ 
    // check for failures on the fly
    // raise exception once a failure condition gets valid
    << saturate agenda with priorities >>
    return true // no accessible states got forbidden!
  catch fail0 fail1 fail2 then
    return false

```

---

## A. Implementation

We present more details on a concrete implementation of the inclusion test for factorized tree automata of Section 4. The same implementation can be used for tree automata without factorization, after conversion into factorized automata. We use a pseudo-functional programming language with imperative state, and have used Objective CAML for implementation in practice.

For simplicity, we restrict ourselves to an inclusion test without dynamic addition of new automata rules. It is not difficult, however, to make the same algorithm incremental in that respect, by returning the complete data structures at the end of the computation, rather than a Boolean value only.

The algorithm applies function `inclusion(A,F)` in Figure 25 to a productive stepwise tree automaton  $A$  with  $\epsilon$ -rules and a deterministic factorized tree automaton  $F$ . It computes the least fixed point of  $D_2(A,F)$  by saturation. The two failure conditions `fail0` and `fail2` are covered by saturation rules (`fail0/a`) and (`fail2/a`). The other failure condition `fail1` is tested either by saturation rule (`fail1/a`), or by using early detection as argued in Section 4.5. Once a failure is detected, an exception is raised in order to exit the saturation loop.

As stated in Section 4.5,  $l_i(p)$  counts the number of literals `f.frbic(p,Q)` inferred so far, and counter  $l_i(p,q)$  the number of occurrences of  $q \in Q$  in literals `f.frbic(p,Q)` seen so far. Their implementations `li(p)` and `li(p,q)` in Figure 26 are counter objects, whose values are initialized to 0. A counter object `C` provides functions `C.val` returning the current value and `C.incr` in-

Figure 26: `<< create counters >>`


---

```

// create counters with initial value 0
forall p ∈ sta(A) do
  l1(p) = counter.new(0)
  l2(p) = counter.new(0)
  forall q ∈ sta(F) do
    l1(p,q) = counter.new(0)
    l2(p,q) = counter.new(0)
// check membership of f.frbi literals to the least fixed point
fun Counters.test(lit)
  case lit
  of f.frb1(p,q) then
    return l1(p).val() > l1(p,q).val()
  of f.frb2(p,q) then
    return l2(p).val() > l2(p,q).val()

```

---

crementing the current value by 1. Function `Counters.test` checks whether its argument  $A, F \models f.frb_i(p, q)$  holds, by comparing the current values of the counters  $l_i(p).val() > l_i(p, q).val()$ . Here, we use a Boolean valued function `>` for comparing integers.

Object `Literals` in Figure 27 collects all literals inferred so far, with the exception of  $frb_i^c$  literals. Function `Literals.mem` tests membership of its argument to this collection. Procedure `Literals.add` adds a literal if it is not yet present in the collection and applies all clauses to it. For literals  $acc(p, q)$ , one first checks whether `fail1` should be raised, either because  $frb(p, q)$  has been inferred before or due to some implied  $f.frb_i$  literal (early detection of failure `fail1`). Second, `fail2` is checked according to  $(fail_2)$ . Third, the literal is put to the collection and onto the agenda with low priority (see below). The addition of literals  $frb(p, q)$  is successful if  $acc(p, q)$  has not been added before. Otherwise exception `fail1` is raised. New literals  $f.acc$  are added to the collection and the agenda with high priority. New literals  $acc(p, \_)$  are treated the same way, except that they are put to the agenda with low priority.

Predicates  $f.frb_i^c$  are always put to the agenda with high priority without further checking. Note that, by this, we permit the addition of the same literal several times. The only operation with those literals will be to increment counters. Also, no rule for those predicates implies the inference of other literals to the fixed point. The halt of the saturation process described below only depends on `acc` and  $f.acc$  predicates. Since such literals are added at most once, halt is guaranteed.

We use an object `Agenda` defined in Figure 28 that stores all literals to which some clauses remain to be applied. Literals in the agenda are either tagged with priority `high` or `low`, as required for early detection of failure



Figure 27: `<< create literal collection >>`

---

```

let heap = Set.new( $\emptyset$ ) in // initialize collection of literals
fun Literals.mem(lit) // test membership of literal to collection
    return heap.mem(lit)
proc Literals.add(lit) // add literal to collection
    case lit
    of acc(p,q) then
        if not Literals.mem(acc(p,q)) then
            if Literals.mem(frb(p,q)) then
                raise fail1 // apply (fail1/a)
            if exists r  $\in$  sta(F) such that q  $\xrightarrow{\epsilon}_F$  r then
                if Counters.test(f.frbsort(q)(p,q))
                    or Counters.test(f.frbsort(r)(p,r)) then
                        raise fail1 // early failure detection
                if p  $\in$  fin(A) and q  $\notin$  fin(F) and r  $\notin$  fin(F) then
                    raise fail2 // apply (fail2/a)
                else //  $\nexists$  r  $\in$  sta(F) such that q  $\xrightarrow{\epsilon}_F$  r
                    ... // same as above but without r
                heap.put(acc(p,q))
                Agenda.put_low(acc(p,q))
            of f.acc(p,q) then
                if not Literals.mem(f.acc(p,q)) then
                    heap.put(f.acc(p,q))
                    Agenda.put_high(f.acc(p,q))
            of acc(p,_) then
                if not Literals.mem(acc(p,_)) then
                    heap.put(acc(p,_))
                    Agenda.put_low(acc(p,_))
            of frb(p,q) then
                if not Literals.mem(frb(p,q)) then
                    if Literals.mem(acc(p,q)) then
                        raise fail1 // apply (fail1/a)
                    heap.put(frb(p,q))
                    Agenda.put_low(frb(p,q))
            of f.frbic(p,Q) then
                Agenda.put_high(f.frbic(p,Q))

```

---

Figure 28: `<< create agenda with priorities >>`


---

```

let high = Stack.new(0) in // define a higher priority stack
let low = Stack.new(0) in // define a lower priority stack
// interface
fun Agenda.nonempty()
    return high.nonempty() and low.nonempty()
fun Agenda.nonempty_high()
    return high.nonempty()
fun Agenda.nonempty_low()
    return low.nonempty()
fun Agenda.get_high()
    return high.pop()
fun Agenda.get_low()
    return low.pop()
proc Agenda.put_high(literal)
    high.push(literal)
proc Agenda.put_low(literal)
    low.push(literal)

```

---

Figure 29: `<< saturate agenda with priorities >>`


---

```

// schedule literals for constant rules
forall a,p such that a → p ∈ rul(A) do
    if exists q ∈ sta(F) such that a → q ∈ rul(F) then
        Literals.add(acc(p,q)) // apply (acc1a)
    else
        raise fail0 // apply (fail0/a)
// saturate agenda with priorities
while Agenda.nonempty() do
    while Agenda.nonempty_high() do
        << apply rules with higher priority >>
    if Agenda.nonempty_low() then
        << apply rules with lower priority >>

```

---

`fail1`. High priority literals may serve in clauses that produce literals of high priority only. The first addition of a literal to the agenda is always with `high`. Once all consequences of high priority are produced, the tag is changed to `low`. Here, we optimize the previous processing by noticing that low priority suffices for treating `acc` and `frb` literals, as well as `f.frbic` literals only have high priorities. The functions of object `Agenda` are `nonempty()`, `get_high()`, `get_low()`, `put_high(lit)`, and `put_low(lit)`.

After initializations, the algorithm starts the saturation process of Figure 29. In a first time, it applies the rule (`acc1a`) for adding accessible states of  $A \times F$  from constants. In parallel it tests for failure `fail0` and raises the appropriate exception if the rule (`fail0/a`) can be applied. During this process, the agenda is filled with low priority `acc` literals.

Figure 30:  $\langle\langle$  apply rules with higher priority  $\rangle\rangle$

---

```

case Agenda.get_high()
of f.acc(p,q) then
  case sort(q)
  of 1 then
    let p1 = p in
      forall p2 such that Af=p1@p2 do
        Literals.add(f.frb2c(p2,Q2f(q))) // apply (f.frb2c)
  of 2 then
    ... // apply (f.frb1c) symmetrically to 1
  Agenda.put_low(f.acc(p,q)) // schedule for low priority operations
of f.frbic(p,Q) then
  li(p).incr()
  forall q ∈ Q do
    li(p,q).incr()

```

---

In a second step, the saturation procedure loops on the agenda and applies the priority policy as follows. It first applies all the rules for literals with the highest priority until it is empty. Then it applies the rules for a unique literal with low priority. This step can indeed involve the scheduling of tasks with a higher priority. The loop stops whenever the whole agenda is empty, or if some exception is raised during saturation.

Applying rules with high priority is described in Figure 30. It does not concern `acc` literals as previously stated. For `f.acc` literals, the highest priority is to infer `f.frbic` literals with respect to sorts  $i \in \{1, 2\}$  by applying (`f.frb1c`) and (`f.frb2c`) rules of Figure 15. Then the literal is scheduled for low priority operations. For literals `f.frbic(p, Q)`, counter  $l_i(\mathbf{p})$  and counters  $l_i(\mathbf{p}, \mathbf{q})$  have to be incremented for all  $q \in Q$  (see Lemma 24). This way, multiple occurrences of such literals in the least fixed point are properly taken into account, as previously noticed.

Applying rules with low priority is described in Figure 31. It concerns only literals `acc` and `f.acc`. For both, the job consists in verifying the necessary conditions to apply all the rules where they appear in the tail of a Datalog clause, namely (`f.acc`), (`acc/2a`), (`acc/4`), (`frb/1a`) and (`frb/2a`) for `acc` literals, and only (`acc/3a`) for `f.acc` since other rules are applied with high priority.

Figure 31:  $\langle\langle$  apply rules with lower priority  $\rangle\rangle$

---

```

case Agenda.get_low()
of acc(p,q) then
  forall r such that  $q \xrightarrow{F}^{\leq 1} r$  do
    Literals.add(f.acc(p,r)) // apply (f.acc)
  forall p' such that  $p \xrightarrow{A} p'$  do
    Literals.add(acc(p',q)) // apply (acc/2a)
    Literals.add(acc(p,-)) // apply (acc/4)
of acc(p,-) then
  // apply (frb/2a)
  let p1 = p in
    forall p2 such that  $A \models p_1 @ p_2$  do
      forall q2  $\notin R_2^F$  do
        Literals.add(frb(p2,q2))
      // apply (frb/1a)
      ... // symmetric to (frb/2a)
of f.acc(p,q) then
  case sort(q)
  of 1 then
    let (p1,q1) = (p,q) in
      forall p2,p' such that  $p_1 @ p_2 \rightarrow p' \in \text{rul}(A)$  do
        forall q2,q' such that  $q_1 @ q_2 \rightarrow q' \in \text{rul}(F)$  do
          if Literals.mem(f.acc(p2,q2)) then
            Literals.add(acc(p',q')) // apply (acc3/a)
  of 2 then
    ... // symmetric to 1

```

---