

Efficient Inclusion Checking for Deterministic Tree Automata and DTDs

Jérôme Champavère, Rémi Gilleron, Aurélien Lemay, and Joachim Niehren

INRIA Futurs and Lille University, LIFL, Mostrare project

Abstract. We present a new algorithm for testing language inclusion $L(A) \subseteq L(B)$ between tree automata in time $O(|A| * |B|)$ where B is deterministic. We extend this algorithm for testing inclusion between automata for unranked trees A and deterministic DTDs D in time $O(|A| * |\Sigma| * |D|)$. No previous algorithms with these complexities exist.

1 Introduction

Language inclusion for tree automata is a basic decision problem that is closely related to universality and equivalence [5,14,15]. Tree automata algorithms are generally relevant for XML document processing [11,17,7,13]. Regarding inclusion checking, a typical application is inverse type checking for tree transducers [10]. Another one is schema-guided query induction [4], 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, not necessarily A . Without this assumption the problem becomes DEXPTIME complete [15]. 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 two deterministic automata $L(A) = L(B)$. The converse might be false, i.e., we cannot rely on polynomial time equivalence tests, as for instance, by comparing number of solutions [14] or minimal deterministic tree automata.

In the case of standard tree automata for ranked trees, the well-known naive algorithm for inclusion 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 nonempty language. The problematic step is to complete B before complementing its final states, since completion might require to add rules for all possible left-hand sides. The overall running time may thus become $O(|A| * |\Sigma| * |B|^n)$, which is exponential in the maximal rank n of function symbols in the signature Σ .

It seems folklore that one can bound the maximal arity of a signature to 2. This can be done by transforming ranked trees into binary trees, and then converting automata for ranked trees into automata for binary trees correspondingly. The problem is to preserve determinism in such a construction, while the size of

automata should grow at most linearly. We show how to obtain such a transformation by stepwise tree automata [3,5]. Thereby we obtain an inclusion test in time $O(|A| * |\Sigma| * |B|^2)$. This is still too much in practice with DTDs, where A and B may be of size 500 and Σ of size 100.

Our first contribution is a more efficient algorithm for standard tree automata on binary trees that verifies inclusion in time $O(|A| * |B|)$ if B is deterministic. This bound is independent of the size of the signature, even though Σ is not fixed. As a second contribution, we show how to test inclusion between automata A for unranked trees and deterministic DTDs D in time $O(|A| * |\Sigma| * |D|)$. Determinism is required by the XML standards. Our algorithm first computes the Glushkov automata of all regular expressions of D in time $O(|\Sigma| * |D|)$, which is possible for deterministic DTDs [2]. The second step is more tedious. We would like to transform the collection of Glushkov automata to a deterministic stepwise tree automaton of the same size. Unfortunately, this seems difficult to achieve, since the usual construction of [9] eliminates ϵ -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 factorized tree automata, which use ϵ -transitions for representing deterministic automata more compactly. We show how to adapt our inclusion test to factorized tree automata and thus to DTDs.

Related Work and Outline. Heuristic algorithms for inclusion between non-deterministic schemas that avoid the high worst-case complexity were proposed in [16]. The complexity of inclusion for various fragments of DTDs and extended DTDs was studied in [8]. The algorithms presented there assume the same types of regular expressions on both sides of the inclusion test. When applied to deterministic DTDs, the same complexity results may be obtainable. Our algorithm permits richer left-hand sides.

We reduce inclusion for the ranked case to the binary case in Section 2. The efficient algorithm for binary tree automata is given in Section 3. In Section 4, we introduce deterministic factorized tree automata and lift the algorithm for inclusion testing to them. Finally in Section 5, we apply them to testing inclusion of automata for unranked trees in deterministic DTDs.

2 Standard Tree Automata for Ranked Trees

We reduce the inclusion problem of tree automata for ranked trees [5] to the case of binary trees with a single binary function symbol.

A ranked signature Σ is a finite set of function symbols $f \in \Sigma$, each of which has an arity $n \geq 0$. 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 of arity n and n trees $t_1, \dots, t_n \in T_\Sigma$.

A *tree automaton* (possibly with ϵ -rules) A over Σ consists of a finite set $s(A)$ of states, a subset $\text{final}(A) \subseteq s(A)$ of final states, and a set $\text{rules}(A)$ of rules of the form $f(p_1, \dots, p_n) \rightarrow p$ or $p' \xrightarrow{\epsilon} p$ where $f \in \Sigma$ has arity n and $p_1, \dots, p_n, p, p' \in s(A)$. We write $p' \xrightarrow{\epsilon}_A p$ iff $p' \xrightarrow{\epsilon} p \in \text{rules}(A)$, $\xrightarrow{\epsilon}_A^*$ for the

reflexive transitive closure of $\xrightarrow{\epsilon}_A$, and $\xrightarrow{\epsilon \leq 1}_A$ for the union of $\xrightarrow{\epsilon}_A$ and the identity relation on $\mathfrak{s}(A)$.

The *size* $|A|$ of A is the sum of the cardinality of $\mathfrak{s}(A)$ and the number of symbols in $\text{rules}(A)$, i.e., $\sum_{f(p_1, \dots, p_n) \rightarrow p \in \text{rules}(A)} (n + 2)$. The cardinality of the signature can be ignored, since our algorithms will not take unused function symbols into account. Every tree automaton A defines an evaluator $\text{eval}_A : T_{\Sigma \cup \mathfrak{s}(A)} \rightarrow 2^{\mathfrak{s}(A)}$ such that $\text{eval}_A(f(t_1, \dots, t_n)) = \{p \mid p_1 \in \text{eval}_A(t_1), \dots, p_n \in \text{eval}_A(t_n), f(p_1, \dots, p_n) \rightarrow p' \in \text{rules}(A), p' \xrightarrow{\epsilon^*_A} p\}$ and $\text{eval}_A(p) = \{p\}$. A tree $t \in T_\Sigma$ is *accepted* by A if $\text{final}(A) \cap \text{eval}_A(t) \neq \emptyset$. The *language* $L(A)$ is the set of trees accepted by A .

A tree automaton is (bottom-up) *deterministic* if it has no ϵ -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.

We will study the inclusion problem for tree automata, whose input consists of a ranked signature Σ , tree automata A with ϵ -rules and deterministic B , both over Σ . Its output is the truth value of $L(A) \subseteq L(B)$. We can deal with this problem by restriction to so-called stepwise signatures $\Sigma_{@}$, consisting of a single binary function symbol $@$ and a finite set of constants $a \in \Sigma$. A stepwise tree automaton [3] is a tree automaton over a stepwise signature.

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

We first encode ranked trees into binary trees via Curryng. Given a ranked signature Σ we define the corresponding signature $\Sigma_{@} = \{@\} \uplus \Sigma$ whereby all symbols of Σ become constants. Curryng is defined by a function $\text{curry} : T_\Sigma \rightarrow T_{\Sigma_{@}}$ which for all trees $t_1, \dots, t_n \in T_\Sigma$ and $f \in \Sigma$ satisfies:

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

For instance, $f(a, g(a, b), c)$ is mapped to $f@a@(g@a@b)@c$ which is infix notation with left-most parenthesis for the tree $@(@(@(f, a), @(@(g, a), b)), c)$. Now we encode tree automata A over Σ into stepwise tree automata $\text{step}(A)$ over $\Sigma_{@}$, such that the language is preserved up to Curryng, i.e., such that $L(\text{step}(A)) = \text{curry}(L(A))$. The states of $\text{step}(A)$ are the prefixes of left-hand sides of rules in A , i.e., words in $\Sigma(\mathfrak{s}(A))^*$:

$$\mathfrak{s}(\text{step}(A)) = \{f q_1 \dots q_i \mid f(q_1, \dots, q_n) \rightarrow q \in \text{rules}(A), 0 \leq i \leq n\} \uplus \mathfrak{s}(A)$$

Its rules extend prefixes step by step by states q_i according to the rule of A . Since constants do not need to be extended, we distinguish two cases in Fig. 1.

Lemma 1. *The encoding of tree automata A over Σ into stepwise tree automata $\text{step}(A)$ over $\Sigma_{@}$ preserves determinism, the tree language modulo Curryng, 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 $\text{step}(B)$.

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

Fig. 1. Transforming ranked tree automata into stepwise tree automata.

3 Stepwise Tree Automata for Binary Trees

We present a new inclusion test that applies to stepwise tree automata over binary trees. We first characterize inclusion into deterministic tree automata, second, express the characterization in Datalog [6] and third, turn it into an efficient algorithm. While the two first steps are easy, the last step is nontrivial.

Characterization of Inclusion. We call a state $p \in s(A)$ *accessible* if there exists a tree t such that $p \in \text{eval}_A(t)$. We call p *co-accessible* if there exists a tree $t \in T_{\Sigma \cup \{p\}}$ with a unique occurrence of p such that $\text{eval}_A(t) \cap \text{final}(A) \neq \emptyset$. A tree automaton is *productive* if all its states are accessible and co-accessible. We denote the *product* of two automata A and B with the same signature by $A \times B$. The state set of $A \times B$ is $s(A) \times s(B)$. For inferring its rules, we assume that B does not have ϵ -rules:

$$\begin{array}{c}
\frac{a \rightarrow p \in \text{rules}(A) \quad a \rightarrow q \in \text{rules}(B)}{a \rightarrow (p, q)} \quad
\frac{p_1 @ p_2 \rightarrow p \in \text{rules}(A) \quad q_1 @ q_2 \rightarrow q \in \text{rules}(B)}{(p_1, q_1) @ (p_2, q_2) \rightarrow (p, q)} \quad
\frac{p' \xrightarrow{\epsilon} p \in \text{rules}(A) \quad q \in s(B)}{(p', q) \xrightarrow{\epsilon} (p, q)}
\end{array}$$

We do not care about final states of $A \times B$ since these are useless in our characterization of inclusion.

Proposition 2. *Inclusion $L(A) \subseteq L(B)$ for a productive stepwise tree automaton A with ϵ -rules and a deterministic stepwise tree automaton B fails iff:*

- fail₀:** *there exists a rule $a \rightarrow p \in \text{rules}(A)$ but no state $q \in s(B)$ such that $a \rightarrow q \in \text{rules}(B)$, or*
- fail₁:** *there exist accessible states (p_1, q_1) and (p_2, q_2) of $A \times B$ and a rule $p_1 @ p_2 \rightarrow p \in \text{rules}(A)$ but no state $q \in s(B)$ such that $q_1 @ q_2 \rightarrow q \in \text{rules}(B)$, or*
- fail₂:** *some accessible state (p, q) of $A \times B$ satisfies $p \in \text{final}(A)$ but $q \notin \text{final}(B)$.*

Proof. If one of the failure conditions holds, then failure of inclusion follows from the hypotheses that A is productive and B deterministic.

For the converse, let us consider a tree t such that $t \in L(A)$ and $t \notin L(B)$. 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 $\text{eval}_B(a) = \emptyset$, thus **fail₀** holds. 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 $p_1 @ p_2 \rightarrow p \in \text{rules}(A)$, since $t \in L(A)$. Since t' is defined as a

$$\begin{array}{c}
(\text{acc}_{/1}) \frac{a \rightarrow p \in \text{rules}(A) \quad a \rightarrow q \in \text{rules}(B)}{\text{acc}(p, q).} \quad (\text{acc}_{/2}) \frac{p' \xrightarrow{\epsilon}_A p \in \quad q \in \text{s}(B)}{\text{acc}(p, q) :- \text{acc}(p', q).} \\
(\text{acc}_{/3}) \frac{p_1 @ p_2 \rightarrow p \in \text{rules}(A) \quad q_1 @ q_2 \rightarrow q \in \text{rules}(B)}{\text{acc}(p, q) :- \text{acc}(p_1, q_1), \text{acc}(p_2, q_2).} \\
(\text{frb}) \frac{p_1 @ p_2 \rightarrow p \in \text{rules}(A) \quad \nexists q. q_1 @ q_2 \rightarrow q \in \text{rules}(B)}{\text{frb}(p_2, q_2) :- \text{acc}(p_1, q_1). \\ \text{frb}(p_1, q_1) :- \text{acc}(p_2, q_2).} \\
(\text{fail}_0) \frac{a \rightarrow p \in \text{rules}(A) \quad \nexists q. a \rightarrow q \in \text{rules}(B)}{\text{fail}_0.} \\
(\text{fail}_1) \frac{p \in \text{s}(A) \quad q \in \text{s}(B)}{\text{fail}_1 :- \text{acc}(p, q), \text{frb}(p, q).} \quad (\text{fail}_2) \frac{p \in \text{final}(A) \quad q \notin \text{final}(B)}{\text{fail}_2 :- \text{acc}(p, q).}
\end{array}$$

Fig. 2. Transforming tree automata A and B into a Datalog program $D_1(A, B)$.

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 $q_1 @ q_2 \rightarrow q \in \text{rules}(B)$. This leads to **fail**₁.
(ii) If $\text{eval}_B(t) \neq \emptyset$ then there exists $q \in \text{eval}_B(t)$; B being deterministic this q is necessarily unique. Since $t \notin L(B)$, $q \notin \text{final}(B)$. Moreover, since $t \in L(A)$, there exists $p \in \text{eval}_A(t) \cap \text{final}(A)$. This leads to **fail**₂. \square

Testing the Characterization. The following efficiency theorem for ground Datalog will be fundamental to all what follows. Given a Datalog program P (without negation), we write $\text{lfp}(P)$ for its least fixed point semantics.

Theorem 1 (Efficiency of Ground Datalog [6]). *For every ground Datalog program P , the least fixed point semantics $\text{lfp}(P)$ can be computed in linear time $O(|P|)$ where the size $|P|$ is the number of symbols in P .*

This result holds even without any bound on the arity of the relation symbols of P , which will be very useful later on. If relation symbols of higher arities are used, the number of their arguments is accounted for by the size of P .

Fig. 2 presents a Datalog program $D_1(A, B)$ that verifies the characterization of $L(A) \subseteq L(B)$ in Proposition 2. Transformation rules $(\text{acc}_{/1})$, $(\text{acc}_{/2})$, and $(\text{acc}_{/3})$ define clauses for accessibility in $A \times B$ through predicate **acc**. The clauses produced by transformation rule (frb) define *forbidden states* of $A \times B$ through predicate **frb**. These are states that lead to **fail**₁ when accessed. Transformation rules (fail_0) , (fail_1) , and (fail_2) define clauses for failures. The characterization of inclusion from Proposition 2 is captured in the following sense:

Proposition 3. *Let A and B be stepwise tree automata for binary trees. If A is productive and B deterministic then:*

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

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

Fig. 3. Grouping (frb) transformations.

The sum of the sizes of the clauses defined by transformation rules ($\text{acc}_{/1}$), ($\text{acc}_{/2}$), ($\text{acc}_{/3}$), (fail_0), (fail_1), and (fail_2) is $O(|A| * |B|)$. The sizes of the clauses defined by transformation rule (frb) sum up to $O(|A| * |\mathfrak{s}(B)|^2)$. The overall size of the ground Datalog program $D_1(A, B)$ is $O(|A| * (|B| + |\mathfrak{s}(B)|^2))$, which may be $O(|A| * |B|^2)$ in the worst case. Therefore, using Theorem 1, inclusion can be decided in time $O(|A| * |B|^2)$.

Efficient Algorithm. This running time is not better than that of the naive algorithm. The square factor is due to the computation of forbidden states for capturing fail_1 . Since (frb) rules cannot be inferred efficiently enough with a Datalog program, we introduce a new predicate frb^c that will group (frb) rules. Using an appropriate data structure, the frb predicates can be induced efficiently from frb^c . The semantics of the latter is given below:

$$A, B \models \text{frb}^c(p, Q) \Leftrightarrow \forall q \in \mathfrak{s}(B) \setminus Q, A, B \models \text{frb}(p, q)$$

Formally, we impose an order $<$ on $\mathfrak{s}(B)$ and consider $\text{frb}^c(p, \{q_1, \dots, q_n\})$ as $(n + 1)$ -ary literals $\text{frb}^c(p, q_{i_1}, \dots, q_{i_n})$ such that $\{q_{i_1}, \dots, q_{i_n}\} = \{q_1, \dots, q_n\}$ and $q_{i_1} < \dots < q_{i_n}$.

In Fig. 3, we propose two transformations ($\text{frb}_{/1}^c$) and ($\text{frb}_{/2}^c$) for inferring frb^c clauses, both of which group (frb) transformations. Note that for every state p , there may be several sets Q such that $\text{frb}^c(p, Q)$ gets inferred. Therefore, we will have to test efficiently whether a state belongs to the union of complements of those state sets. This will be further detailed.

Let us consider the transformation of tree automata A and B into a ground Datalog program $D_2(A, B)$ defined by transformation rules ($\text{acc}_{/1}$), ($\text{acc}_{/2}$), ($\text{acc}_{/3}$), ($\text{frb}_{/1}^c$), ($\text{frb}_{/2}^c$), (fail_0), and (fail_2). The clauses producing acc , fail_0 , and fail_2 of $D_1(A, B)$ and $D_2(A, B)$ are identical and their number is in $O(|A| * |B|)$. The number of frb^c clauses introduced by rule ($\text{frb}_{/1}^c$) is in $O(|A| * |\mathfrak{s}(B)|)$ but the size of each such clause is $n + 1$ which in the worst case could be $|\mathfrak{s}(B)| + 1$, and symmetrically for ($\text{frb}_{/2}^c$). The overall size of all frb^c clauses, however, is bounded by the overall number of acc clauses, which in turn is bounded by $O(|A| * |B|)$, too! To see this, we can rewrite ($\text{frb}_{/2}^c$) as shown in Fig. 4, such that the corresponding ($\text{acc}_{/3}$) clauses are inferred simultaneously (and these don't overlap). Therefore the overall size of $D_2(A, B)$ is in $O(|A| * |B|)$.

$$\begin{array}{c}
\left. \begin{array}{l}
q_1 @ q_2^1 \rightarrow q^1 \in \text{rules}(B) \\
\vdots \\
q_1 @ q_2^n \rightarrow q^n \in \text{rules}(B)
\end{array} \right\} \text{all the rules for } q_1 \\
\hline
p_1 @ p_2 \rightarrow p \in \text{rules}(A) \\
\text{frb}^c(p_2, \{q_2^1, \dots, q_2^n\}) :- \text{acc}(p_1, q_1). \\
\text{acc}(p, q^1) :- \text{acc}(p_1, q_1), \text{acc}(p_2, q_2^1). \\
\vdots \\
\text{acc}(p, q^n) :- \text{acc}(p_1, q_1), \text{acc}(p_2, q_2^n).
\end{array}$$

Fig. 4. Rewriting grouping rules for complexity analysis of (frb^c_2) clauses.

Inclusion Test. First it computes $\text{lfp}(D_2(A, B))$ in time $O(|A| * |B|)$. If fail_0 or fail_2 belong to $\text{lfp}(D_2(A, B))$ then inclusion does not hold, so **false** is returned. Otherwise, we test for fail_1 in a second step, by checking for all states $\text{acc}(p, q) \in \text{lfp}(D_2(A, B))$ whether there is a $\text{frb}^c(p, Q) \in \text{lfp}(D_2(A, B))$ such that $q \in \mathfrak{s}(B) \setminus Q$. If so, **false** is returned, otherwise **true**.

We have to prove that the second step can be done in time $O(|A| * |B|)$. For every state $p \in \mathfrak{s}(A)$, there are some state sets Q_1, \dots, Q_m such that for $1 \leq j \leq m$, $\text{frb}^c(p, Q_j) \in \text{lfp}(D_2(A, B))$ and we have to check efficiently whether some state q is in $\bigcup_{j=1}^m \mathfrak{s}(B) \setminus Q_j$. For this, we define a data structure $\text{bad_states}(p)$ as an array T of size $|\mathfrak{s}(B)| + 1$. Counters are indexed by elements in $\mathfrak{s}(B)$ and one counter is indexed by 0. All counter values are set to 0 initially. The initialization of all $\text{bad_states}(p)$ can be done in $O(|\mathfrak{s}(A)| * |\mathfrak{s}(B)|)$. For every p and every Q_j such that $\text{frb}^c(p, Q_j) \in \text{lfp}(D_2(A, B))$, counter $T[0]$ is incremented by 1 and counter $T[q]$ is incremented by 1 for all $q \in Q_j$, which can be done in time $O(|Q_j|)$. As the overall size of frb^c clauses in $\text{lfp}(D_2(A, B))$ is in $O(|A| * |B|)$, the computation of all $\text{bad_states}(p)$ can be done in time $O(|A| * |B|)$.

It remains to test for all states $\text{acc}(p, q) \in \text{lfp}(D_2(A, B))$ whether there exists a $\text{frb}^c(p, Q) \in \text{lfp}(D_2(A, B))$ such that $q \in \mathfrak{s}(B) \setminus Q$. This is done by checking whether $T[q] < T[0]$ in $\text{bad_states}(p)$ (see, e.g., Fig. 5). Indeed, if $\{Q \mid \text{frb}^c(p, Q) \in \text{lfp}(D_2(A, B))\} = \{Q_1, \dots, Q_m\}$ then $\text{bad_states}(p)$ is defined such that $T[q] = T[0]$ iff $q \in \bigcap_{j=1}^m Q_j$, thus $T[q] < T[0]$ iff $q \in \bigcup_{j=1}^m \mathfrak{s}(B) \setminus Q_j$. Each such test costs $O(1)$ so the overall time is bounded by $O(|A| * |B|)$.

This concludes the inclusion test for stepwise tree automata, for productive A and deterministic B . Every tree automaton can be made productive in linear time. Higher arities can be reduced to 2 by Proposition 1. This yields:

	0	q ₁	q ₂	q ₃	q ₄	q ₅	...
p	2	0	1	2	1	0	...

Fig. 5. Data structure bad_state .

Here, it has been set up $\text{frb}^c(p, \{q_2, q_3\})$ and $\text{frb}^c(p, \{q_3, q_4\})$. When $\text{frb}^c(p, \{q_2, q_3\})$ is set up, $\text{bad_state}(p)(0)$, $\text{bad_state}(p)(q_2)$ and $\text{bad_state}(p)(q_3)$ are incremented. We have for instance $\text{frb}(p, q_4)$ since $\text{bad_state}(p)(q_4) = 1$ is lower than $\text{bad_state}(p)(0) = 2$. In fact, the only *not forbidden* state is (p, q_3) because q_3 belongs to the intersection of $\{q_2, q_3\}$ and $\{q_3, q_4\}$.

Theorem 2. *Let A and B be standard tree automata for ranked trees of some signature Σ possibly with ϵ -rules. If B is deterministic, inclusion $L(A) \subseteq L(B)$ can be decided in time $O(|A| * |B|)$ independently of the size of Σ .*

4 Factorized Tree Automata

We next relax the determinism assumption on B in a controlled manner, that will be crucial to deal with DTDs. We replace B by deterministic factorized automata, that we introduce. These are automata with ϵ -rules, that represent deterministic automata in a compact manner.

Definition 1. *A factorized tree automaton F over a stepwise signature Σ is a stepwise tree automaton with ϵ -rules and a partition $\mathfrak{s}(F) = \mathfrak{s}_1(F) \uplus \mathfrak{s}_2(F)$ such that if $q_1 @ q_2 \rightarrow q$ in $\text{rules}(F)$ then $q_1 \in \mathfrak{s}_1(F)$ and $q_2 \in \mathfrak{s}_2(F)$.*

We say that q is of sort i in F if $q \in \mathfrak{s}_i(F)$. 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 $\text{ta}(F)$ without ϵ -rules that recognizes the same language. Both automata have the same signature and states; the rules of $\text{ta}(F)$ are inferred as follows from those of F :

$$(E_1) \frac{a \rightarrow q \in \text{rules}(F)}{a \rightarrow q \in \text{rules}(\text{ta}(F))} \quad (E_2) \frac{q_1 \xrightarrow{F}^* r_1 \quad q_2 \xrightarrow{F}^* r_2 \quad r_1 @ r_2 \rightarrow q \in \text{rules}(F)}{q_1 @ q_2 \rightarrow q \in \text{rules}(\text{ta}(F))}$$

We set $\text{final}(\text{ta}(F)) = \{q \mid q \xrightarrow{F}^* r, r \in \text{final}(F)\}$. Note that the size of $\text{ta}(F)$ may be $O(|\text{rules}(F)| * |\mathfrak{s}(F)|^2)$ 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.

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

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

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

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

Proposition 4. *The tree automaton $\text{ta}(F)$ represented by a deterministic factorized tree automaton F is deterministic.*

Proof. Let $B = \text{ta}(F)$ which by construction is free of ϵ -rules. For every constant $a \in \Sigma$, the uniqueness of q such that $a \rightarrow q \in \text{rules}(B)$ follows from \mathbf{d}_0 . For every $q_1 @ q_2 \rightarrow q$ in $\text{rules}(B)$ 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 $q_1 \xrightarrow{F}^* r_1$ at most one r_2 of sort 2 such that $q_2 \xrightarrow{F}^* r_2$. Condition \mathbf{d}_0 implies that there exists at most one state q such that $r_1 @ r_2 \rightarrow q \in \text{rules}(B)$. \square

$$\begin{array}{c}
(\text{acc}/_{3a}) \frac{p_1 @ p_2 \rightarrow p \in \text{rules}(A) \quad q_1 @ q_2 \rightarrow q \in \text{rules}(F)}{\text{acc}(p, q) :- \text{f.acc}(p_1, q_1), \text{f.acc}(p_2, q_2)}. \quad (\text{f.acc}) \frac{p \in \mathfrak{s}(A) \quad q \xrightarrow{\epsilon}_{F}^{\leq 1} r}{\text{f.acc}(p, r) :- \text{acc}(p, q)}. \\
(\text{f.frb}_2^c) \frac{p_1 @ p_2 \rightarrow p \in \text{rules}(A) \quad q_1 \in \mathfrak{s}_1(F)}{\text{f.frb}_2^c(p_2, Q_2^F(q_1)) :- \text{f.acc}(p_1, q_1)}. \quad (\text{f.frb}_1^c) \frac{p_1 @ p_2 \rightarrow p \in \text{rules}(A) \quad q_2 \in \mathfrak{s}_2(F)}{\text{f.frb}_1^c(p_1, Q_1^F(q_2)) :- \text{f.acc}(p_2, q_2)}. \\
(\text{frb}_2^c) \frac{p_1 @ p_2 \rightarrow p \in \text{rules}(A) \quad q_1 \in \mathfrak{s}(F)}{\text{frb}_2^c(p_2, R_2^F) :- \text{acc}(p_1, q_1)}. \quad (\text{frb}_1^c) \frac{p_1 @ p_2 \rightarrow p \in \text{rules}(A) \quad q_2 \in \mathfrak{s}(F)}{\text{frb}_1^c(p_1, R_1^F) :- \text{acc}(p_2, q_2)}. \\
(\text{frb}_{/1a}^c) \frac{p_1 @ p_2 \rightarrow p \in \text{rules}(A) \quad q_2 \notin R_2^F}{\text{frb}^c(p_1, \emptyset) :- \text{acc}(p_2, q_2)}. \quad (\text{frb}_{/2a}^c) \frac{p_1 @ p_2 \rightarrow p \in \text{rules}(A) \quad q_1 \notin R_1^F}{\text{frb}^c(p_2, \emptyset) :- \text{acc}(p_1, q_1)}. \\
(\text{fail}_{2a}) \frac{p \in \text{final}(A) \quad \forall r. q \xrightarrow{\epsilon}_{F}^{\leq 1} r \Rightarrow r \notin \text{final}(F)}{\text{fail}_2 :- \text{acc}(p, q)}.
\end{array}$$

The clauses from $(\text{acc}/_1)$, $(\text{acc}/_2)$ and (fail_0) in $D_2(A, F)$ belong to $D_3(A, F)$, too. We use sets of states $Q_2^F(q_1) = \{q' \mid q_1 @ q' \rightarrow q'' \in \text{rules}(F)\}$, $Q_1^F(q_2)$ symmetrically, and sets of states reaching a sort $R_i^F = \{q \mid \exists r \in \mathfrak{s}_i(F). q \xrightarrow{\epsilon}_{F}^{\leq 1} r\}$.

Fig. 6. Inferring clauses of Datalog program $D_3(A, F)$ simulating $D_2(A, B)$.

We fix a stepwise tree automaton A and a deterministic factorized tree automaton F , and let us $B = \text{ta}(F)$. We now show how to test language inclusion $L(A) \subseteq L(B)$ from A and F without computing B . This is done by the ground Datalog program $D_3(A, F)$ of Fig. 6.

We need new predicates for properties of F in order to infer corresponding properties of B . The accessibility predicate f.acc for F subsumes the accessibility predicate acc for B . Subsumption may be proper as stated by the rule (f.acc) of $D_3(A, F)$. *Vice versa*, we infer accessibility in F from accessibility in B according to the rule $(\text{acc}/_{3a})$. Rules $(\text{acc}/_1)$ and $(\text{acc}/_2)$ of $D_2(A, F)$ remain valid for accessibility in B , too.

Lemma 2. $\text{acc}(p, q) \in \text{lfp}(D_3(A, F))$ iff $\text{acc}(p, q) \in \text{lfp}(D_2(A, B))$

We need to refine predicate frb into predicates frb_1 and frb_2 that take sorts into account, and corresponding predicates f.frb_1 and f.frb_2 in the factorized case. Their semantics can be defined as follows, where A, B are tree automata and F is a factorized tree automaton.

$$\begin{array}{l}
A, B \models \text{frb}_2(p_2, q_2) \Leftrightarrow \exists p, p_1, q_1. A, B \models \text{acc}(p_1, q_1), p_1 @ p_2 \rightarrow p \in \text{rules}(A), q_2 \notin Q_2^B(q_1) \\
A, F \models \text{f.frb}_2(p_2, r_2) \Leftrightarrow \exists p, p_1, r_1. A, F \models \text{f.acc}(p_1, r_1), p_1 @ p_2 \rightarrow p \in \text{rules}(A), r_2 \notin Q_2^F(r_1)
\end{array}$$

The semantics of frb_1 and f.frb_1 are symmetric. The relation to the previous predicate frb is that $A, B \models \text{frb}(p, q)$ if and only if $A, B \models \text{frb}_1(p, q) \vee \text{frb}_2(p, q)$.

The Datalog program $D_3(A, F)$ infers for sorts $i \in \{1, 2\}$ literals with predicates f.frb_i^c that are to be understood by grouping of f.frb_i literals, and similarly frb_i^c by grouping of frb_i . These grouping mechanisms account for sorts, since complementation is with respect to sorts.

$$\begin{aligned}
A, F \models \text{f.frb}_i^c(p, Q) &\Leftrightarrow \forall q \in \mathfrak{s}_i(F) \setminus Q. A, F \models \text{f.frb}_i(p, q) \\
A, F \models \text{frb}_i^c(p, Q) &\Leftrightarrow \forall q \in \mathfrak{s}_i(F) \setminus Q. A, F \models \text{frb}_i(p, q)
\end{aligned}$$

The clauses produced by (f.frb_i^c) and (frb_i^c) are sound with respect to this semantics for deterministic F 's. This is easier to see for (f.frb_i^c) than for (frb_i^c) . We prove it by the next lemma. For states p, q, r and sorts $i \in \{1, 2\}$ we define:

$$\begin{aligned}
A, B \vdash \text{frb}_i(p, q) &\text{ iff } \exists Q \subseteq \mathfrak{s}(B) \setminus \{q\}. \text{frb}^c(p, Q) \in \text{lfp}(D_2(A, B)) \text{ via } (\text{frb}_{/i}^c) \\
A, F \vdash \text{f.frb}_i(p, r) &\text{ iff } \exists R \subseteq \mathfrak{s}(F) \setminus \{r\}. \text{f.frb}_i^c(p, R) \in \text{lfp}(D_3(A, F)) \\
A, F \vdash \text{frb}_i(p, r) &\text{ iff } \exists R \subseteq \mathfrak{s}(F) \setminus \{r\}. \text{frb}_i^c(p, R) \in \text{lfp}(D_3(A, F)) \\
&\text{ or } \text{frb}^c(p, \emptyset) \in \text{lfp}(D_3(A, F)) \text{ via } (\text{frb}_{/ia}^c)
\end{aligned}$$

Lemma 3 (Core). $A, B \vdash \text{frb}_i(p, q)$ iff $A, F \vdash \text{frb}_i(p, q)$ or the unique state r of sort i with $q \xrightarrow{F}^{\leq 1} r$ exists and satisfies $A, F \vdash \text{f.frb}_i(p, r)$.

Since the size of $D_3(A, F)$ is in $O(|A| * |F|)$ we can compute the set of all $\text{f.frb}_1^c(p, R)$ and $\text{f.frb}_2^c(p, R)$ literals in $\text{lfp}(D_3(A, F))$ in time $O(|A| * |F|)$. It remains to infer the induced literals $\text{f.frb}_i(p, r)$ literals in an efficient manner. Computing all $A, F \vdash \text{f.frb}_i(p, r)$ from $\text{lfp}(D_3(A, F))$ can be done by the same clever algorithm as before for deducing all $A, B \vdash \text{frb}(p, q)$ given $\text{lfp}(D_2(A, B))$. Note, however, that we now need two different data structures for the two sorts.

Theorem 3. For a stepwise tree automaton with ϵ -rules A and a deterministic factorized tree automaton F over the same signature, inclusion $L(A) \subseteq L(F)$ can be decided in time $O(|A| * |F|)$.

5 Automata for Unranked Trees and DTDs

We lift our results to deterministic tree automata for unranked trees possibly with factorization, so that they become applicable to deterministic DTDs.

An unranked signature Σ is a finite set of symbols (without arity restrictions). The set T_Σ^u of unranked trees over Σ is the least set that contains all pairs $a(t_1, \dots, t_n)$ where $a \in \Sigma$ and (t_1, \dots, t_n) is a possibly empty sequence of unranked trees in T_Σ^u . Currying carries over literally from ranked to unranked trees. This yields a bijection $\text{curry} : T_\Sigma^u \rightarrow T_{\Sigma_\otimes}$. Thus, we can reuse stepwise tree automata to recognize languages of unranked trees, and as before, we can factorize them. So, as a corollary of Theorem 3 we have:

Corollary 2. For a stepwise tree automaton for unranked trees A and a deterministic factorized tree automaton for unranked trees F over the same signature Σ , $L(A) \subseteq L(F)$ can be decided in time $O(|A| * |F|)$ independently of $|\Sigma|$.

Note that hedge automata [5] can be translated in linear time to stepwise tree automata with ϵ -rules [9]. The automata of [12] support factorization, too. We finally show how to convert deterministic DTDs D to deterministic factorized tree automata in time $O(|\Sigma| * |D|)$, so that we can reuse our algorithm for testing inclusion in deterministic DTDs. Factorization avoids the quadratic blowup from translating hedge to stepwise automata [9].

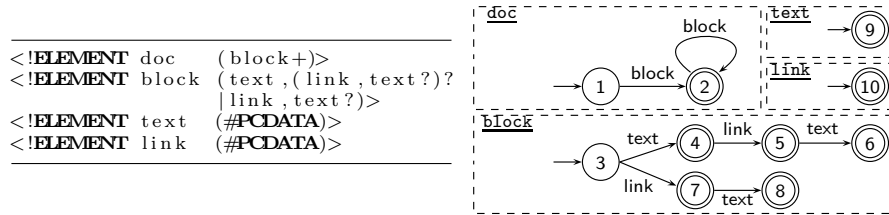


Fig. 7. An example DTD and the corresponding Glushkov automata.

A DTD D with elements in a set Σ is a function mapping letters $a \in \Sigma$ to regular expressions e over Σ , what we write $a \rightarrow_D e$ in this case. One of these elements is the distinguished start symbol. 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 . An example in XML syntax is given in Fig. 7. The set of elements of D is $\Sigma = \{\text{doc}, \text{block}, \text{text}, \text{link}\}$, of which the element 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 [2]. This is equivalent to say that all corresponding Glushkov automata are deterministic, which is a requirement by the W3C. See Fig. 7 for the example.

Theorem 4 (Brüggemann-Klein [1]). *The collection of Glushkov automata for a deterministic DTD D over Σ can be computed in time $O(|\Sigma| * |D|)$.*

We transform a 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 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., $q \xrightarrow{\epsilon} a \in \text{rules}(F)$. The only final state of F is the start symbol of the DTD D . The result is a

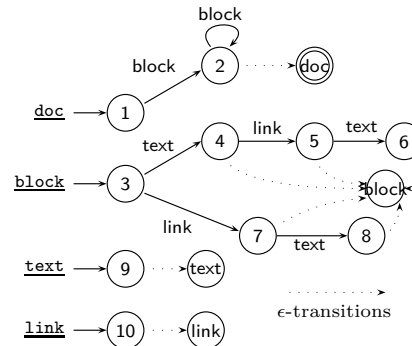


Fig. 8. The deterministic factorized tree automaton for the DTD in Fig. 7.

factorized tree automaton, as for instance in Fig. 8. This needs time of at most $O(|\Sigma| * |D|)$. Note that the size of the example automaton would grow quadratically, when eliminating ϵ -edges. For every $a \in \Sigma$, there is $a \rightarrow q \in \text{rules}(F)$ 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 $q@a \rightarrow q' \in \text{rules}(F)$.

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

Theorem 5. *Deterministic DTDs D over Σ can be translated to deterministic factorized tree automata that recognize the same language in time $O(|\Sigma| * |D|)$.*

Corollary 3. *Language inclusion of hedge automata A over Σ in deterministic DTDs D with elements in Σ can be decided in time $O(|A| * |\Sigma| * |D|)$.*

References

1. A. Brüggemann-Klein. Regular expressions to finite automata. *Theoretical Computer Science*, 120(2):197–213, 1993.
2. A. Brüggemann-Klein and D. Wood. One-unambiguous regular languages. *Information and Computation*, 142(2):182–206, 1998.
3. J. Carme, J. Niehren, and M. Tommasi. Querying unranked trees with stepwise tree automata. In 19th RTA, vol. 3091 of LNCS, 105 – 118. 2004.
4. J. Champavère, R. Gilleron, A. Lemay, and J. Niehren. Towards schema-guided XML query induction. In *ICML CAGI Workshop*, 2007.
5. H. Comon et al. Tree automata techniques and applications. Available on: <http://tata.gforge.inria.fr>, 2007.
6. E. Dantsin, T. Eiter, G. Gottlob, and A. Voronkov. Complexity and expressive power of logic programming. *ACM computing surveys*, 33(3):374–425, 2001.
7. S. Maneth, A. Berlea, T. Perst, and H. Seidl. XML type checking with macro tree transducers. In *24th PODS*, pages 283–294, 2005.
8. W. Martens, F. Neven, and T. Schwentick. Complexity and decision problems for XML Schemas and chain regular expressions. 2008. Journal extension of MFCS'04.
9. W. Martens and Joachim Niehren. On the minimization of XML schemas and tree automata for unranked trees. *J. of Comp. and Sys. Sci.*, 73(4):550–583, 2007.
10. T. Milo, D. Suciu, and V. Vianu. Type checking XML transformers. *J. of Comp. and Sys. Sci.*, 1(66):66–97, 2003.
11. A. Neumann and H. Seidl. Locating matches of tree patterns in forests. In FSTTCS, pages 134–145, 1998.
12. Stefan Raeymaekers. *Information Extraction from Web Pages Based on Tree Automata Induction*. PhD thesis, Katholieke Universiteit Leuven, 2008.
13. T. Schwentick. Automata for XML—a survey. *J. of Comp. and Sys. Sci.*, 73(3):289–315, 2007.
14. H. Seidl. Deciding equivalence of finite tree automata. *SIAM Journal on Computing*, 19(3):424–437, 1990.
15. H. Seidl. Haskell overloading is DEXPTIME-complete. *Information Processing Letters*, 52(2):57–60, 1994.
16. A. Tozawa and M. Hagiya. XML schema containment checking based on semi-implicit techniques. In *Int. Conf. on Impl. and Appl. of Automata*, 2003.
17. S. Dal Zilio and D. Lugiez. XML schema, tree logic and sheaves automata. In 18th RTA, volume 2706 of LNCS, pages 246–263. 2003.