



HAL
open science

Automata for Analyzing and Querying Compressed Documents

Barbara Fila, Siva Anantharaman

► **To cite this version:**

Barbara Fila, Siva Anantharaman. Automata for Analyzing and Querying Compressed Documents. 2006. hal-00088776v1

HAL Id: hal-00088776

<https://hal.science/hal-00088776v1>

Preprint submitted on 4 Aug 2006 (v1), last revised 15 Dec 2006 (v3)

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.



UNIVERSITE D'ORLEANS

Faculté des Sciences

LIFO

Laboratoire d'Informatique Fondamentale d'Orléans
4, rue Léonard de Vinci, BP 6759
F-45067 Orléans Cedex 2
FRANCE

Rapport de Recherche

[www : http://www.univ-orleans.fr/SCIENCES/LIFO/](http://www.univ-orleans.fr/SCIENCES/LIFO/)

Automata for Analyzing and Querying Compressed Documents

Barbara FILA, LIFO, Orléans (Fr.)
Siva ANANTHARAMAN, LIFO, Orléans (Fr.)

Rapport N° **2006-03**

Automata for Analyzing and Querying Compressed Documents

Barbara Fila, Siva Anantharaman

LIFO - Université d'Orléans (France),
e-mail: {fila, siva}@univ-orleans.fr

Abstract. In a first part of this work, tree/dag automata are defined as extensions of (unranked) tree automata which can run indifferently on trees or dags; they can thus serve as tools for analyzing or querying any semi-structured document, whether or not given in a compressed format. In a second part of the work, we present a method for evaluating positive unary queries, expressed in terms of Core XPath axes, on any dag t representing an XML document possibly given in a compressed form; the evaluation is done directly on t , without unfolding it into a tree. To each Core XPath query of a certain basic type, we associate a word automaton; these automata run on the graph of dependency between the non-terminals of the minimal straightline regular tree grammar associated to the given dag t , or along complete sibling chains in this grammar. Any given positive Core XPath query can be decomposed into queries of the basic type, and the answer to the query, on the dag t , can then be expressed as a sub-dag of t whose nodes are suitably labeled under the runs of such automata.

Keywords: Tree automata, Tree grammars, Dags, XML, Core XPath.

1 Introduction

Several algorithms have been optimized in the past, by using structures over dags instead of over trees. Tree automata are widely used for querying XML documents (e.g., [8, 9, 15, 16]); on the other hand, the notion of a compressed XML document has been introduced in [2, 7, 12], and a possible advantage of using dag structures for the manipulation of such documents has been brought out in [12]. It is legitimate then to investigate the possibility of using automata over dags instead of over trees, for querying compressed XML documents.

Dag automata (DA) were first introduced and studied in [5]; a DA was defined there as a natural extension of tree automaton, i.e. as a bottom-up tree automaton running on dags; and the language of a DA was defined as the set of dags that get accepted under (bottom-up) runs, defined in the usual sense; the emptiness problem for DAs was shown there to be NP-complete, and the membership problem proved to be in NP; but the problem of stability under complementation of the class of dag automata –closely linked with that of determinization– was left open. These two issues have since been settled negatively in [1]: the reason is that the set of all terms (trees) represented by the set of dags accepted by a *non-deterministic* DA is not necessarily a regular tree language; a consequence is that the class of tree languages recognized by DAs (as sets of accepted dags) is a *strict* superclass of the class of regular tree languages. It is well-known however, that answers to MSO-definable queries on (semi-)structured trees form regular tree languages ([18]); it is thus necessary to define the languages of DAs in a manner *different* from that of [5, 1], if they are to serve as tools for analyzing and querying a document, independently of whether it is given in a (partially or fully) compressed format, or as a tree. Our first aim in this work is therefore to redefine the notion of the language of a DA suitably, with such an objective.

For achieving that, we first present (in Section 2) the notion of a compressed document as a *tree/dag* (*trdag*, for short), designating a directed acyclic graph that may be partially or fully compressed. The terminology *trdag* has been chosen to distinguish it from that of *tdag* employed in [1]; this latter term will be employed in this paper when referring to a fully compressed dag. A Tree/Dag automaton (TDA, for short) is then defined as an automaton which runs on trdags. The essential differences with the DAs of [1] are the following: (i) our TDAs can be unranked, and (ii) although the transition rules of a TDA look quite like those of the DAs in [1], or those of TAs, a run of a TDA on any given *trdag* t will carry with it not only assignments of states to the nodes of t , but also to the edges of t ; runs will be so defined that a TDA accepts any given *trdag* t if and only if it accepts the tree \hat{t} obtained by uncompressing t , as a tree automaton running on the tree \hat{t} , in the usual sense.

In the second part of the paper, we present an approach based on word automata for evaluating queries on trdags that represent XML documents in a partially or fully compressed format; the terms ‘trdag’ and ‘document’ will therefore be considered synonymous in the sequel. Any given *trdag* t is first seen as equivalent to a minimal straightline regular tree grammar \mathcal{L}_t , that one can naturally associate with t , cf. e.g., [3, 4]. From the grammar \mathcal{L}_t , we construct the graph of dependency \mathcal{D}_t between its non-terminals, and also the *chiblings* (linear graphs formed of complete chains of sibling non-terminals) of \mathcal{L}_t . The word automata that we build below will run on \mathcal{D}_t or the chiblings of \mathcal{L}_t , rather than on the document t itself.

We shall only consider *positive* unary queries expressed in terms of Core XPath axes. (The view we adopt allows us to define the various axes of Core XPath on compressed documents, in a manner which does not modify their semantics on trees.) For evaluating any such query on any document (*trdag*) t , we proceed as follows. We first break up the given query into basic sub-queries of the form $Q = // *[axis::\sigma]$ where *axis* is a Core XPath axis of a certain type. To each such basic query Q , we associate a word automaton \mathbf{A}_Q . The automaton \mathbf{A}_Q runs on the graph \mathcal{D}_t when *axis* is non-sibling, and on the chiblings of \mathcal{L}_t when *axis* is a sibling axis. An essential point in our method is that the runs of \mathbf{A}_Q are guided by some well-defined semantics for the nodes traversed, indicating whether the current node answers Q , or is on a path leading to some other node answering Q . The automaton, though not deterministic, is made effectively unambiguous by defining a suitable priority relation between its transitions, based on the semantics. A basic query Q can then be evaluated in one single top-down pass of \mathbf{A}_Q , under such an unambiguous run. An arbitrary positive unary Core XPath query can be evaluated on t by combining the answers to its various basic sub-queries, and its answer set is expressed as a sub-*trdag* of t , whose nodes get labeled in conformity with the semantics. It is important to note that the evaluation is performed on the *given* *trdag* t ; as such, on two different trdags corresponding to two different compressions of a same XML tree, the answers obtained may *not* be the same, in general.

The paper is structured as follows: Section 2 presents the notions of trdags, and of Tree/Dag automata. In Section 3, we construct from any *trdag* t its normalized straightline regular tree grammar \mathcal{L}_t , as well as the dependency graph \mathcal{D}_t and the chiblings of \mathcal{L}_t ; these will be seen as rooted labeled acyclic graphs (*rlags*, for short); the basic notions of Core XPath are also recalled. Section 4 is devoted to the construction of the word automata for any basic Core XPath query, based on the semantics, and an illustrative example. In Section 5 we prove that the runs of these automata, uniquely and effectively determined under a maximal priority condition, generate the answers to the queries. Section 6 shows how a non basic (composite, or imbricated) Core XPath query can be evaluated in a stepwise fashion. In the appendices, we show how to translate the

‘usual’ Core XPath queries into one in ‘standard’ form, on which our approach is applicable; we also present a polynomial time algorithm for constructing the maximal priority run for any basic query automaton over any given document (trdag), with a complexity bound of degree 3 on the number of nodes of the trdag; a complete illustrative example, on a composite imbricated query, is given.

2 Tree/Dag Automata

Definition 1 A tree/dag (trdag for short) over a not necessarily ranked alphabet Σ is a rooted dag (directed acyclic graph) $t = (Nodes(t), Edges(t))$, where, for any node $u \in Nodes(t)$:

- u has a name $name_t(u) = name(u) \in \Sigma$;
- the edges going out of any node are ordered;
- and if $name(u)$ is ranked, then the number of outgoing edges at u is the rank of $name(u)$.

Given any node u on a trdag t , the notion of the sub-trdag of t rooted at u is defined as usual, and denoted as $t|_u$. If v is any node, $\gamma(v) = u_1 \dots u_n$ will denote the string of all its not necessarily distinct children nodes; for every $1 \leq i \leq n$, the i -th outgoing edge from v to its i -th child node $u_i \in \gamma(v)$ will be denoted as $e(v, i)$; we shall also write then $v \xrightarrow{i} u_i$; the set of all outgoing (resp. incoming) edges at any node v will be denoted as Out_v (resp. In_v); and for any node u other than the root node of t , we set: $Parents(u) = \{v \in Nodes(t) \mid u \text{ is a child of } v\}$.

A trdag t will be said to be a tree iff for every node u on t other than the root, $Parents(u)$ is a singleton. For any trdag t , we define the set $Pos(t)$ as the set of all the positions $pos_t(u)$ of all its nodes u , these being defined recursively, as follows: if u is the root node on t , then $pos_t(u) = \epsilon$, otherwise, $pos_t(u) = \{\alpha.i \mid \alpha \in pos_t(v), v \text{ is a parent of } u, u \text{ is an } i\text{-th child of } v\}$. The elements of $Pos(t)$ are words over natural integers,

The function $name_t$ is extended naturally to the positions in $Pos(t)$ as follows: for every $u \in Nodes(t)$ and $\alpha \in pos_t(u)$, we set $name_t(\alpha) = name_t(u)$. Given a trdag t , we define its tree-equivalent as a tree \hat{t} such that: $Pos(\hat{t}) = Pos(t)$, and for every $\alpha \in Pos(t)$ we have $name_t(\alpha) = name_{\hat{t}}(\alpha)$. A trdag is said to be a tdag, or *fully compressed*, iff for any two different nodes u, u' on t , the two sub-dags $t|_u$ and $t|_{u'}$ have non-isomorphic tree-equivalents; otherwise, the trdag is said to be *partially compressed*. For example, the tree to the left of Figure 1 is the tree-equivalent of the partially compressed trdag to the right, and also to the fully compressed tdag to the middle.

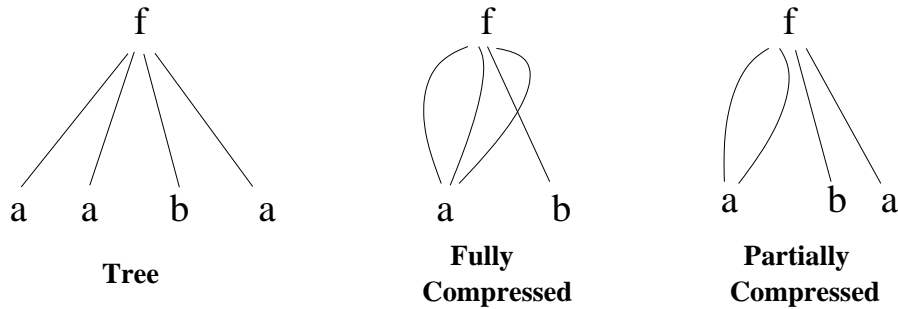


Fig. 1. tree, tdag, and trdag

We define now the notion of a Tree/Dag automaton, first over a ranked alphabet Σ , to facilitate understanding. The definition is then easily extended to the unranked case.

Definition 2 A *Tree/Dag automaton (TDA, for short)* over a ranked alphabet Σ is a tuple (Σ, Q, F, Δ) , where Q is a finite non-empty set of states, $F \subseteq Q$ is the set of final (or accepting) states, and Δ is a set of transition rules of the form: $f(q_1, \dots, q_k) \rightarrow q$, where $f \in \Sigma$ is of rank k , and $q_1, \dots, q_k, q \in Q$.

It will be convenient to write the transition rules of a TDA in a different (but equivalent) form: a transition of the form $f(q_1, \dots, q_k) \rightarrow q$ is also written as $(f, q_1 \dots q_k) \rightarrow q$, where $q_1 \dots q_k$ is seen as a word in Q^* , of length = $rank(f)$ in the ranked case. The notion of a TDA is then extended easily to the unranked case, where the signature symbols naming the nodes are not assumed to be of fixed rank: it suffices to define the transitions to be of the form $(f, \omega) \rightarrow q$, where $\omega \in Q^*$; we may assume wlog that ω is a ‘*’-regular expression on Q not involving ‘+’, by replacing a rule $(f, \omega + \omega') \rightarrow q$, by the two rules $(f, \omega) \rightarrow q, (f, \omega') \rightarrow q$.

A TDA is said to be bottom-up *deterministic* iff whenever there are two transition rules of the form $(f, \omega) \rightarrow q, (f, \omega') \rightarrow q'$, with $q \neq q'$, we have necessarily $\omega \cap \omega' = \emptyset$; otherwise it is said to be *non-deterministic*. We also agree to denote the transitions of the form $(f, \emptyset) \rightarrow q$ simply as $f \rightarrow q$, and refer to them as *initial* transitions.

For defining the notion of runs of TDAs on a trdag in a bottom-up style, we need some preliminaries. Let \mathbf{A} be a TDA with state set Q and transition set Δ . Suppose t is a trdag and assume given a map $M : Edges(t) \rightarrow Q$. If u is any node on t with $u_1 \dots u_n$ as the string of all its (not necessarily distinct) children, the string $M(\mathbf{e}(u, 1)) \dots M(\mathbf{e}(u, n))$, formed of states assigned by M to the outgoing edges at u , will be denoted as $M(Out_u)$. We then define, recursively in a bottom-up style, a binary relation at u on the states of Q , with respect to (w.r.t. or wrt, for short) the given map M ; this relation, denoted as $\triangleleft_u^M = \triangleleft_u$, is defined as follows:

Definition 3 Let \mathbf{A}, t, M be as above, and u any given node on the trdag t .

- If u is a leaf with $name(u) = a$, then $q \triangleleft_u q'$ iff whenever $a \rightarrow q \in \Delta$ we also have $a \rightarrow q' \in \Delta$;
- otherwise $q \triangleleft_u q'$ iff:
 - (i) $(name(u), M(Out_u)) \rightarrow q$ is an instance of a transition rule in Δ ; i.e., Δ has a rule $(name(u), \omega) \rightarrow q$ such that $M(Out_u)$ is in ω ;
 - (ii) there exists a map $\sigma_{q'} : Q \rightarrow Q$, such that:
 - $\sigma_{q'}(q) = q'$, and the rule $(name(u), \sigma_{q'}(M(Out_u))) \rightarrow q'$ is also an instance of a transition rule in Δ ;
 - for any edge $\mathbf{e} : u \xrightarrow{i} u' \in Out_u$, we have: $M(\mathbf{e}) \triangleleft_{u'} \sigma_{q'}(M(\mathbf{e}))$.

Definition 4 Let $\mathbf{A} = (\Sigma, Q, F, \Delta)$ be any given TDA, and t any given trdag. A run of \mathbf{A} on t is a pair (r, M) , where $r : Nodes(t) \rightarrow Q$ and $M : Edges(t) \rightarrow Q$ are maps such that the following conditions hold, at any node u on t :

(1) if $name(u) = f$, then the rule $(f, M(Out_u)) \rightarrow r(u)$ is an instance of a transition rule in Δ ;

(2) there is an incoming edge $\mathbf{e} \in In_u$ with $M(\mathbf{e}) = r(u)$; and for every $\mathbf{e}' \in In_u$ such that $M(\mathbf{e}') = q' \neq q = r(u)$, we have $q \triangleleft_u^M q'$

A run (r, M) is accepting on a trdag t iff $r(\epsilon) \in F$, i.e., r maps the root-node of t to an accepting state. A trdag t is accepted by a TDA iff there is an accepting run on t . The language of a TDA is the set of all trdags that it accepts.

Remark 1. i) Note that if t is a tree, then In_u is singleton at every non-root node u on t , so a run (r, M) of any TDA on t can be identified with its first component r ; we get then the usual notion of runs of tree automata on trees.

Example 1. Over the unranked signature $\{a, f, g\}$ consider a TDA \mathbf{A} , with the following transitions:

$$\begin{aligned}
a \rightarrow p, \quad b \rightarrow q', \quad b \rightarrow p, \quad b \rightarrow q, \\
(a, p) \rightarrow q, \quad (a, q) \rightarrow p, \\
(g, qQ^*) \rightarrow q, \quad (g, pq) \rightarrow p, \\
(f, qpq) \rightarrow q_{fin}, \quad (f, pQ^*) \rightarrow q_{fin},
\end{aligned}$$

with $Q = \{p, q, q', q_{fin}\}$, and q_{fin} as the unique accepting state. An accepting bottom-up run of \mathbf{A} on a tdag is depicted on the left of Figure 2, and on its right, the “same” run as seen on the tree equivalent of the tdag.

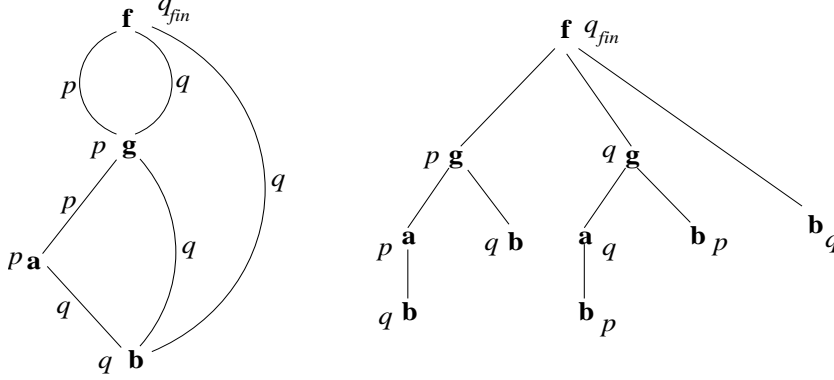


Fig. 2. A bottom-up accepting run of the TDA of Example 1 on a tdag, and the same seen on its tree equivalent.

A few comments on the above run may be of help: we start with assigning state q to the leaf node b , under r ; the assignments of state q under M to all the incoming edges at this node b poses no problem; we can then assign state p to node a , and subsequently also p to the node g , under r , via the transition rule $(g, pq) \rightarrow p$; we then assign p under M to the first incoming edge at g ; assigning state q under M to the second incoming edge at g is valid since we have: $q \triangleleft_a p$ and $p \triangleleft_b q$; reaching q_{fin} at the root-node is trivial via the last transition rule. (Note that we could have as well assigned p under M to the second incoming edge at g , with no conditions to check, then reach q_{fin} .)

Remark 1 (contd.) ii) Unlike the DAs of [5] or [1], the following bottom-up non-deterministic TDA: $a \rightarrow q_1$, $a \rightarrow q_2$, $f(q_1, q_2) \rightarrow q_a$, with q_0, q_1, q_a as states where q_a is accepting, has a non-empty language: as a TDA it accepts $f(a, a)$.

For a *deterministic* TDA, we have the following result (as expected):

Proposition 1 *Let \mathbf{A} be a bottom-up deterministic TDA, and t any given tdag; then there is at most one run of \mathbf{A} on t .*

Proof. Let Q be the set of states of \mathbf{A} , and $M : Edges(t) \rightarrow Q$ any given map assigning states to the edges on t . We shall show by induction that the hypothesis of determinism on \mathbf{A} implies that, at any node u on t , the binary relation $\triangleleft_u^M = \triangleleft_u$ defined above (Definition 3), w.r.t. the map M , is the identity relation on the set Q . The proposition will then follow from conditions (1) and (2) on runs, cf. Definition 4 (we will get, in particular, that for every incoming edge \mathbf{e} at u , $M(\mathbf{e})$ must be the same as $r(u)$; so the run can be identified with its first component r).

The induction will be on a non-negative integer d_u –that we define at any node u of t , and refer to as its *height* on t – as the maximal number of arcs on t from u to the leaf nodes. If $d_u = 0$, then u is a leaf node; that \triangleleft_u is the identity relation on Q in this case is immediate, from the determinism of \mathbf{A} , and the definition of \triangleleft_u . So, assume that $d_u > 0$, and let $v_1 \dots v_n$ be the string of all the

children nodes of u on t . By the inductive hypothesis, for every $i, 1 \leq i \leq n$, the relation \triangleleft_{v_i} is identity; it follows then, from the conditions (i) and (ii) of the relation \triangleleft_u , that this latter must also be the identity relation on Q . \square

We may now formulate the principal result of the first part of this paper:

Proposition 2 *i) A TDA accepts a trdag t if and only if it accepts the tree equivalent of t .*

ii) The emptiness problem for a TDA is decidable in time P w.r.t. its number of states.

iii) The uniform membership problem for a TDA is decidable in time NP (resp. time P) w.r.t. its number of states, and the number of edges (resp. and the number of positions) on the given trdag.

Proof. Let \hat{t} be the tree equivalent of the trdag t . (It is immediate that \hat{t} is uniquely determined, up to a tree isomorphism.) The ordered bisimulation relation between the sets of nodes of t and \hat{t} can actually be seen as a natural surjective map from $Nodes(\hat{t})$ onto $Nodes(t)$ that can be defined recursively, in a bottom-up style; in what follows, this map will be referred to as the compression map, and denote it as \mathbf{c} .

Property i): For proving the ‘only if’ part, one uses the following reasoning, coupled with induction on the height function at the nodes of t (defined in the proof of the previous proposition): Let (r, M) be an accepting run of the given TDA on the trdag t ; consider a node s on the tree equivalent \hat{t} , of which the node u on t is the image under the compression map \mathbf{c} ; let $r(u) = q$ under the given run of the TDA on t ; then, for every state q' of the TDA such that $q \triangleleft_u^M q'$, one can construct a partial run of the TDA –seen as a usual tree automaton– on the tree \hat{t} , climbing up from a leaf below s on \hat{t} , to the node s (for an illustrative example, see the tree to the right of Figure 2).

For proving the ‘if’ part of Property i), we start with a given accepting run \hat{r} of the given TDA, as a bottom-up tree automaton running in the usual sense on the tree \hat{t} ; from this run \hat{r} , we shall construct a run (r, M) of the TDA on the trdag t as follows: Consider any given node u on t , and let $\{s_1, \dots, s_k\}$ be the set $\mathbf{c}^{-1}(u)$ of all nodes on \hat{t} such that u is the image of each of them under the compression map $\mathbf{c}: Nodes(\hat{t}) \rightarrow Nodes(t)$; the set In_u , of incoming edges at u on the trdag t , contains then exactly k edges, say $\mathbf{e}_1, \dots, \mathbf{e}_k$ in some (arbitrarily chosen) order; we then set, for any $i, 1 \leq i \leq k, M(\mathbf{e}_i) = \hat{r}(s_i)$; and set $r(u)$ as any (arbitrarily chosen) state among the $\hat{r}(s_i), i, 1 \leq i \leq k$. It is then easily checked that the pair of maps (r, M) gives an accepting run of the TDA on the trdag t .

Properties ii) and iii) follow, in the ranked case, from the proof of i) and the results of TATA ([6]), Chapter 1; in the unranked case, one can either employ a reasoning based on reduction to the ranked case as in [10], or appeal directly to the results of [13]. (Note: the number of positions on a trdag is the same as the size of its tree equivalent.) \square

3 Querying Compressed Documents: Preliminaries

Given a trdag t , one can naturally construct a regular tree grammar associated with t , which is *straightline* (cf. [4]), in the sense that there are no cycles on the dependency relations between its non-terminals, and each non-terminal produces exactly one sub-trdag of t . Such a grammar will be denoted as \mathcal{L}_t , if it is *normalized* in the following sense:

(i) for every non-terminal A_i of \mathcal{L}_t , there is exactly one production of the form $A_i \rightarrow f(A_{j_1}, \dots, A_{j_k})$, where $i < j_r$ for every $1 \leq r \leq k$; we shall then set $Sons(A_i) = \{A_{j_1}, \dots, A_{j_k}\}$, and $symb_{\mathcal{L}_t}(A_i) = f$;

(ii) the number of non-terminals is the number of nodes on t .

Such a normalized grammar \mathcal{L}_t is uniquely defined up to a renaming of the non-terminals. For instance, for the trdag t to the left of Figure 3 we get the following normalized grammar:

$$A_1 \rightarrow f(A_2, A_3, A_4, A_5, A_2), \quad A_2 \rightarrow c, \quad A_3 \rightarrow a(A_5), \quad A_4 \rightarrow b, \quad A_5 \rightarrow b.$$

Such a grammar is easily constructed from t , for instance by using a standard algorithm which computes the ‘depth’ of any node (as the maximal distance from the root), to number the non-terminals so as to satisfy condition (i) above.

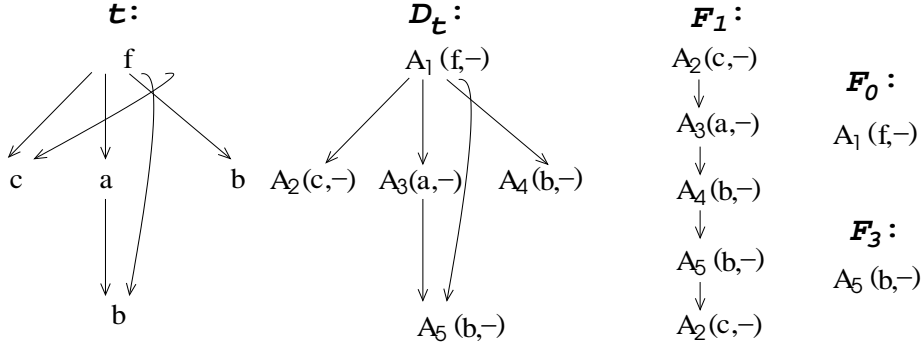


Fig. 3. trdag t , associated rlag \mathcal{D}_t , and chiblings of \mathcal{L}_t

The *dependency graph* of the normalized grammar \mathcal{L}_t associated with t , and denoted as \mathcal{D}_t , consists of nodes named with the non-terminals $A_i, 1 \leq i \leq n$, and *one single* directed arc from any node A_i to a node A_j whenever A_j is a son of A_i . The root of \mathcal{D}_t is by definition the node named A_1 . The notion of *Sons* of the nodes on \mathcal{D}_t is derived in the obvious way from that defined above on \mathcal{L}_t .

Furthermore, to any production $A_i \rightarrow f(A_{j_1}, \dots, A_{j_k})$ of \mathcal{L}_t , we associate a rooted linear graph composed of k nodes respectively named A_{j_1}, \dots, A_{j_k} , with root at A_{j_1} and such that for all $l \in \{2, \dots, k\}$ the node named A_{j_l} is the son of the node named $A_{j_{l-1}}$. This graph will be called the *chibling* of \mathcal{L}_t associated with the (unique) A_i -production; it is denoted as \mathcal{F}_i . We also define a further chibling denoted \mathcal{F}_0 , as the linear graph with a single node named A_1 , where A_1 is the axiom of \mathcal{L}_t .

In the sequel, we designate by \mathcal{G} either \mathcal{D}_t or any of the chiblings \mathcal{F} of \mathcal{L}_t . We complete any of these acyclic graphs \mathcal{G} into a rooted labeled acyclic graph (*rlag*, for short), by attaching to each node u on \mathcal{G} , with $\text{name}(u) = A_i$, a label denoted $\text{label}(u)$, and defined as $\text{label}(u) = (\text{symb}_{\mathcal{L}_t}(A_i), -)$; cf. Figure 3.

3.1 Positive Core XPath Queries on trdags

In this paper we restrict our study to positive Core XPath queries on trdags. Recall that Core XPath is the navigational segment of XPath, and is based on the following axes of XPath (cf. [10, 19]): **self**, **child**, **parent**, **ancestor**, **descendant**, **following-sibling**, **preceding-sibling**. A location expression is defined as a predicate of the form $[\text{axis}::b]$, where **axis** is one of the above axes, and b is a symbol of Σ . Given any trdag t over Σ , a context node u on t and $b \in \Sigma$, the semantics for **axis** is defined by evaluating this predicate at u . The semantics for the axes **self**, **child**, **descendant** are easily defined, exactly as on trees (cf. [19]). For defining the semantics of the remaining axes, we first recall that $\text{Parents}(u) = \{v \in \text{Nodes}(t) \mid u \text{ is a child of } v\}$.

Definition 5 *Given a context node u on a trdag t , and $b \in \Sigma$:*

i) $[\text{parent}::b]$ evaluates to true at u , if and only if there exists a b -named node in $\text{Parents}(u)$;

ii) $[\text{ancestor}::b]$ evaluates to true at u , iff either $[\text{parent}::b]$ evaluates to true at u , or there exists a node $v \in \text{Parents}(u)$ such that $[\text{ancestor}::b]$ evaluates to true at v ;

iii) $[\text{following-sibling}::b]$ evaluates to true at u , iff there exists a b -named node u' , and a node v on t such that $\gamma(v)$ is of the form $\dots u \dots u' \dots$;

iv) $[\text{preceding-sibling}::b]$ evaluates to true at u , iff there exists a b -named node u' , and a node v on t such that $\gamma(v)$ is of the form $\dots u' \dots u \dots$.

For the ‘composite’ axes **descendant-or-self** and **ancestor-or-self**, the semantics are then deduced in an obvious manner. We shall also need position predicates of the form $[\text{position}()=i]$; their semantics is that the expression $[\text{child}::b \text{ } [\text{position}()=i]]$ evaluates to true at a context node u , iff: $[\text{child}::b]$ evaluates to true at u , and u is an i -th child of some parent.

Positive Core XPath query expressions are usually defined in the literature (cf. e.g., [7]), as those generated by the following grammar:

$$\begin{aligned} A &::= \text{self} \mid \text{child} \mid \text{descendant} \mid \text{parent} \mid \text{ancestor} \mid \\ &\quad \text{preceding-sibling} \mid \text{following-sibling} \\ S_{can} &::= A::\sigma \mid \text{position}()=i \mid S_{can} \text{ and } S_{can} \mid S_{can} \text{ or } S_{can} \\ E_{can} &::= A::*[S_{can}] \mid E_{can}[E_{can}] \\ Q_{can} &::= /S_{can} \mid /E_{can} \mid Q_{can}/Q_{can} \end{aligned}$$

We shall refer to the query expressions generated by this grammar as *canonical*; they can be shown to be of the type $/C_1/C_2/\dots/C_n$, where each C_i is of the form $A::\sigma[X_{can}]$, or of the form $A::\sigma[X_{can}] \text{ conn } A'::\sigma'[X'_{can}]$, with $\text{conn} \in \{\text{and}, \text{or}\}$, and $X_{can}, X'_{can} \in \{S_{can}, E_{can}, \text{true}\}$; we agree here to identify $A::\sigma[\text{true}]$ with $A::\sigma$.

Any such positive Core XPath query expression can be translated into one that is in ‘standard form’, i.e., where the format of the sub-queries is of the type ‘ $\text{axis}::b$ ’; we formalize this idea now. We shall refer to the axes **self**, **child**, **descendant**, **parent**, **ancestor**, **preceding-sibling**, **following-sibling** as *basic*. A basic Core XPath query is a query of the form $//*[\text{axis}::\sigma]$, where **axis** is a basic axis. More generally, the queries we propose to evaluate on trdags are defined formally as the expressions Q_{std} generated by the following grammar, where σ stands for any node name on the documents, or for $*$ (meaning ‘any’):

$$\begin{aligned} A &::= \text{self} \mid \text{child} \mid \text{descendant} \mid \text{parent} \mid \text{ancestor} \mid \\ &\quad \text{preceding-sibling} \mid \text{following-sibling} \\ S &::= A::\sigma \mid \text{position}()=i \mid S \text{ and } S \mid S \text{ or } S \mid \text{Root} \\ E &::= A::*[S] \mid E[E] \\ Q_{std} &::= //* \mid //*[S] \mid //*[E] \end{aligned}$$

Core XPath queries Q_{std} of the format generated by this grammar are said to be in *standard form*; to be able to handle any positive Core XPath query with such a grammar, we have introduced a special predicate called **Root**, deemed true only at the root node of the trdag considered.

By the *evaluation* of a given query expression Q on any trdag t , we mean the assignment: $t \mapsto$ the set of all context nodes on t where the expression Q evaluates to true (following the conventions of Definition 2); this latter set is also called the *answer* for Q on t . Two given queries Q_1, Q_2 are said to be *equivalent* iff, on any trdag t , the answer sets for Q_1 and Q_2 are the same. Any positive Core XPath query Q_{can} can be translated into an equivalent one in standard form; e.g., $/c[\text{following-sibling}::g]/d$ is equivalent to $//*[self::d \text{ and } \text{parent}::*[\text{Root and self}::c \text{ } [\text{following-sibling}::g]]]$ in standard form. An inductive procedure performing such a translation in the general case (of linear complexity w.r.t. the number of location steps in Q_{can}) is given in Appendix I. The following proposition results from Definition 5.

Proposition 3 (1) For any set of nodes X on a trdag t , and any axis A , we have: $A(X) =$

$$\bigcup_{\substack{x \in X, \alpha \in \text{post}_t(x) \\ \alpha = i_1 \dots i_k}} \{ / \text{child}:: * [\text{position}() = i_1] / \dots / \text{child}:: * [\text{position}() = i_k] / A:: * \}$$

(2) For any trdag t , and any node with name b on t , we have:

$$(i) // * [\text{preceding}:: b] = \bigcup_u \{ \text{descendant-or-self}(\text{following-sibling}(u), // * [\text{self}:: u \text{ and } (\text{descendant}:: b \text{ or } \text{self}:: b)]) \}$$

$$(ii) // * [\text{following}:: b] = \bigcup_u \{ \text{descendant-or-self}(\text{preceding-sibling}(u), // * [\text{self}:: u \text{ and } (\text{descendant}:: b \text{ or } \text{self}:: b)]) \}$$

Finally, following [2], for any set S of nodes on t , the sets of nodes `following(S)` and `preceding(S)` can now be defined formally, as follows:

$$\text{following}(S) = \text{descendant-or-self}(\text{following-sibling}(\text{ancestor-or-self}(S))),$$

$$\text{preceding}(S) = \text{descendant-or-self}(\text{preceding-sibling}(\text{ancestor-or-self}(S))).$$

Note: Unlike on a tree, the `ancestor`, `descendant`, `following`, `self` and `preceding` axes do *not* partition the set of nodes on a trdag t , in general.

4 Automata for the Basic Core XPath Queries

4.1 The Semantics of the Approach

We first consider basic Core XPath queries. Composite or imbricated queries will subsequently be evaluated in a stepwise fashion; see Section 6.

To any basic query $Q = // * [\text{axis}:: \sigma]$, we shall associate a word automaton (actually a transducer), referred to as \mathbf{A}_Q . It will run top-down, on the rlag \mathcal{D}_t if `axis` is non-sibling, and on each of the chiblings \mathcal{F} of \mathcal{L}_t otherwise. In either case, a run will attach, to any node traversed, a pair of the form (t, x) , where the component t of the pair has the intended semantics of selection or not, by Q , of the corresponding node on t , and the component x will be a 1 or 0, with the intended semantics that $x = 1$ iff the corresponding node on t has a descendant answering Q . At the end of the run, $\text{label}(u)$, at any node u of \mathcal{D}_t , will be replaced by a new label derived from the ll-pairs attached to u by the run.

To formalize these ideas, we introduce a set of new symbols $L = \{s, \eta, \top, \top'\}$ referred to as *llabels* (the term ‘llabel’ is used so as to avoid confusion with the term label). We define *ll-pairs* as elements of the set $L \times \{0, 1\}$, and the states of \mathbf{A}_Q as elements of the set $\{\text{init}\} \cup (L \times \{0, 1\})$. For any Q , the automaton \mathbf{A}_Q is over the alphabet $\Sigma \cup \{s, \eta\}$, has *init* as its initial state, and has no final state. The set Δ_Q of transitions of \mathbf{A}_Q will consist of rules of the form $(q, \tau) \rightarrow q'$ where $q \in \{\text{init}\} \cup (L \times \{0, 1\})$, $q' \in (L \times \{0, 1\})$, and $\tau \in \Sigma \cup \{s, \eta\}$.

For any rlag \mathcal{G} , we define a function $\text{llab}: \text{Nodes}(\mathcal{G}) \rightarrow \Sigma \cup \{s, \eta\}$, by setting $\text{llab}(u) = \pi_1(\text{label}(u))$, the first component of $\text{label}(u)$. The automaton \mathbf{A}_Q associated to a basic query $Q = // * [\text{axis}:: \sigma]$ will run *top-down* on the rlag \mathcal{G} , where \mathcal{G} is \mathcal{D}_t if `axis` is a basic non-sibling axis, and \mathcal{G} is any chibling \mathcal{F} of \mathcal{L}_t if `axis` is a basic sibling axis. A *run* of \mathbf{A}_Q on \mathcal{G} is a map $r: \text{Nodes}(\mathcal{G}) \rightarrow L \times \{0, 1\}$, such that, for every $u \in \text{Nodes}(\mathcal{G})$, the following holds:

- if u is $\text{root}_{\mathcal{G}}$, then the rule $(\text{init}, \text{llab}(u)) \rightarrow r(u)$ is in Δ_Q ;
- otherwise, for every $v \in \gamma(u)$ the rules $(r(u), \text{llab}(v)) \rightarrow r(v)$ are all in Δ_Q .

(Note: when `axis` is non-sibling, this amounts to requiring that, for any node v , the state $r(v)$ must be in conformity with the states $r(u)$ for *every* parent node u of v , with respect to the rules in Δ_Q .)

From the run of the automaton \mathbf{A}_Q and from the states it attaches to the nodes of \mathcal{D}_t , we will deduce, at every node u of t , a well-determined ll-pair as (a new) label at u , via the natural bijection between $Nodes(t)$ and $Nodes(\mathcal{D}_t)$. The ll-pairs thus attached to the nodes of t will have the following semantics (where x stands for the name of the node u on t , corresponding to the ‘current’ node on \mathcal{D}_t):

- $(\top', 1) : x = \sigma$, current node on t is selected by (i.e., is an answer for) Q ;
- $(\top, 1) : x = \sigma$, current node is *not* selected, but has a selected descendant;
- $(\top, 0) : x = \sigma$, current node is *not* selected, and has *no* selected descendant;
- $(s, 1) : x \neq \sigma$, current node is selected;
- $(\eta, 1) : x \neq \sigma$, current node is *not* selected, but has a selected descendant;
- $(\eta, 0) : x \neq \sigma$, current node is *not* selected, and has *no* selected descendant.

Only the nodes on \mathcal{D}_t , to which the run of A_Q associates the labels $(s, 1)$ or $(\top', 1)$, correspond to the nodes of t that will get selected by the query Q . The ll-pairs with boolean component 1 will label the nodes of \mathcal{D}_t corresponding to the nodes of t which are on a path to an answer for the query Q ; thus the automata \mathbf{A}_Q will have *no* transitions from any state with boolean component 0 to a state with boolean component 1. Moreover, with a view to define runs of such automata which are unique (or unambiguous in a sense that will be presently made clear), we define the following *priority* relations between the ll-pairs:

$$(\eta, 0) > (\eta, 1) > (s, 1), \quad \text{and} \quad (\top, 0) > (\top, 1) > (\top', 1).$$

A run of the automaton \mathbf{A}_Q will label any node u on \mathcal{G} with an ll-pair either from the group $\{(\top, 0), (\top, 1), (\top', 1)\}$ or from the group $\{(\eta, 0), (\eta, 1), (s, 1)\}$; and this group is determined by $llab(u)$.

For ease of presentation, we agree to set $\eta' := s$, and often denote either of the above two groups of ll-pairs under the uniform notation $\{(l, 0), (l, 1), (l', 1)\}$, where $l \in \{\eta, \top\}$, with the ordering $(l, 0) > (l, 1) > (l', 1)$.

We shall construct a run r of \mathbf{A}_Q on \mathcal{G} that will be uniquely determined by the following *maximal priority* condition:

(MP): at any node v on \mathcal{G} , $r(v)$ is the maximal ll-pair (t, x) for the ordering $>$ in the group $\{(l, 0), (l, 1), (l', 1)\}$ determined by $llab(v)$, such that \mathbf{A}_Q contains a transition rule of the form $(r(u), llab(v)) \rightarrow (t, x)$, for *every* parent u of v .

Such a run will assign a label with boolean component 1 only to the nodes corresponding to those of the minimal sub-trdag t containing the root of t and all the answers to Q on t .

4.2 Re-labeling of \mathcal{D}_t by the Runs of \mathbf{A}_Q

We first consider a non-sibling basic query Q on a given document t , and a given run r of the automaton \mathbf{A}_Q on the \mathcal{D}_t ; at the end of the run, the nodes on \mathcal{D}_t will get re-labeled with new ll-pairs, computed as below for every $u \in Nodes(\mathcal{D}_t)$:

$$\begin{aligned} lab_r(u) &= (s, 1) \text{ iff } r(u) \in \{(s, 1), (\top', 1)\}, \\ lab_r(u) &= (\eta, 1) \text{ iff } r(u) \in \{(\eta, 1), (\top, 1)\}, \\ lab_r(u) &= (\eta, 0) \text{ iff } r(u) \in \{(\eta, 0), (\top, 0)\}. \end{aligned}$$

The rlag obtained in this manner from \mathcal{D}_t , following the run r and the associated re-labeling function lab_r , will be denoted as $r(\mathcal{D}_t)$.

For a basic query Q over a sibling axis, the situation is a little more complex, because several different nodes on one chibling of \mathcal{L}_t can have the same name (non-terminal), or several different chiblings can have nodes named by the same non-terminal, or both. Thus, to any node of \mathcal{D}_t , named with a non-terminal A , will correspond in general a *set* of ll-pairs, assigned by the various runs of \mathbf{A}_Q to the A -named nodes on the various chiblings of \mathcal{L}_t . We therefore proceed as follows: for every complete set \hat{r} of runs of \mathbf{A}_Q , formed of one run $r_{\mathcal{F}}$ on each

chibling \mathcal{F} , we will define $\widehat{r}(\mathcal{D}_t)$ as the re-labeled rlag derived from \mathcal{D}_t , under \widehat{r} . With that purpose we associate to \widehat{r} and any $u \in \text{Nodes}(\mathcal{D}_t)$, a set of ll-pairs:

$$ll_{\widehat{r}}(u) = \bigcup_{r_{\mathcal{F}} \in \widehat{r}} \{r_{\mathcal{F}}(v) \mid v \in \text{Nodes}(\mathcal{F}), \text{ and } \text{name}(v) = \text{name}(u)\}.$$

We then derive, at each node of \mathcal{D}_t a unique ll-pair in conformity with the semantics of our approach, by using the following function:

$$\begin{aligned} \lambda_{\widehat{r}}(u) = s &\iff ll_{\widehat{r}}(u) \cap \{(s, 1), (\top', 1)\} \neq \emptyset, \\ \lambda_{\widehat{r}}(u) = \eta &\iff ll_{\widehat{r}}(u) \cap \{(s, 1), (\top', 1)\} = \emptyset. \end{aligned}$$

From \mathcal{D}_t and this function $\lambda_{\widehat{r}}$, we next derive an rlag $\lambda_{\widehat{r}}(\mathcal{D}_t)$ by re-labeling each node u on \mathcal{D}_t with the pair $(\lambda_{\widehat{r}}(u), -)$. And finally we define $\widehat{r}(\mathcal{D}_t)$ as the rlag obtained from $\lambda_{\widehat{r}}(\mathcal{D}_t)$, by running on it the automaton for the basic non-sibling query $/**[\text{self}::s]$, as indicated at the beginning of this subsection. In practical terms, such a run amounts in essence to setting, as the second component of $\text{label}(u)$ at any node u , the boolean 1 iff u is on a path to some node with llab s , and 0 otherwise. All these details are illustrated with an example in the following subsection.

4.3 The Automata

We first present the automata for the basic queries $/**[\text{self}::\sigma]$ and for $/**[\text{following-sibling}::\sigma]$, and give an illustrative example using the former for $\sigma = s$, and the latter for $\sigma = b$. The automata for the other basic queries are given after the example.

- Automata: for $/**[\text{self}::\sigma]$ and for $/**[\text{following-sibling}::\sigma]$

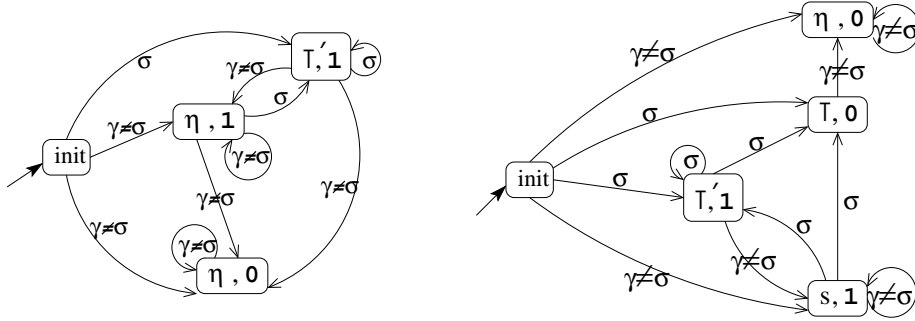


Figure 4 below illustrates the evaluation of $Q = /**[\text{following-sibling}::b]$, on the trdag t of Figure 3. We first use the automaton for the basic query $/**[\text{following-sibling}::\sigma]$ with $\sigma = b$, and then the automaton for $/**[\text{self}::\sigma]$ with $\sigma = s$. The sub-trdag of t , formed of nodes corresponding to those of $\widehat{r}(\mathcal{D}_t)$ with labels having boolean component 1, contains all the answers to Q on t .

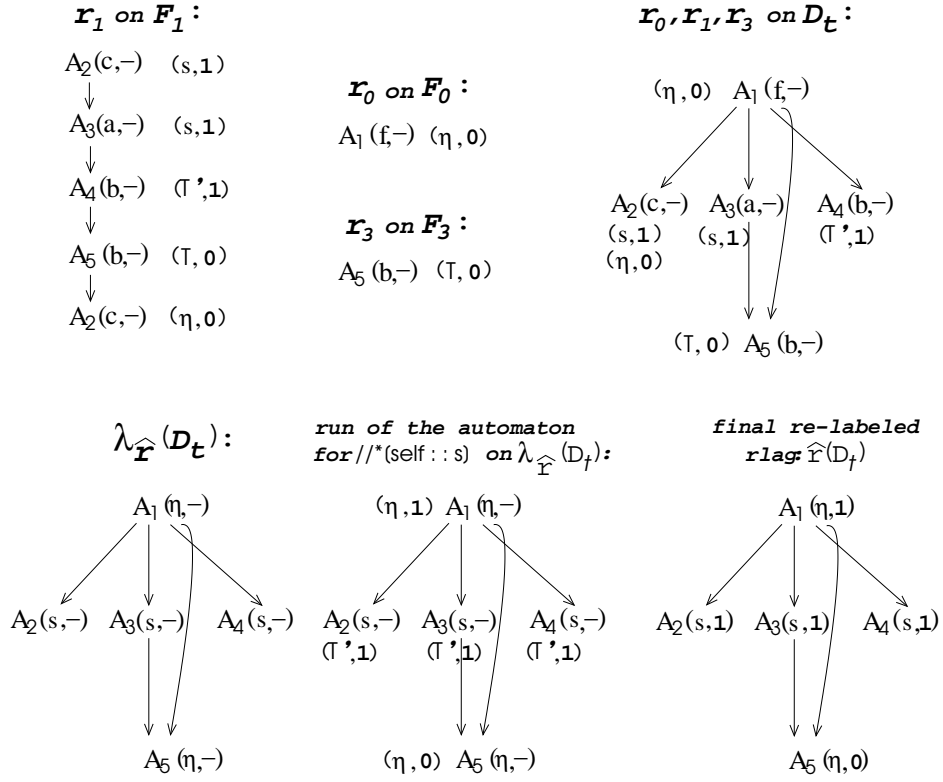
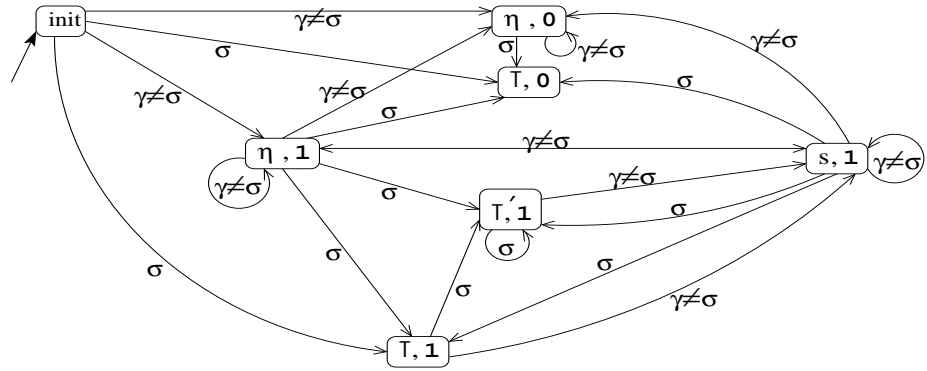
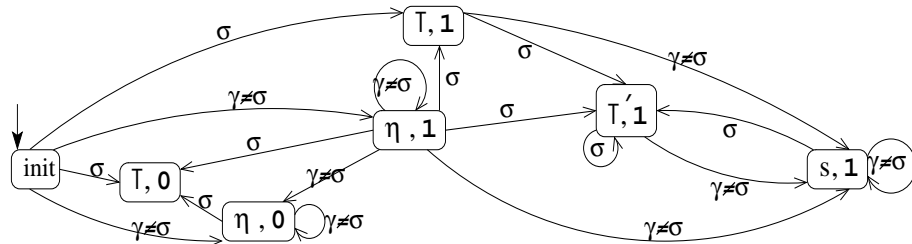


Fig. 4.

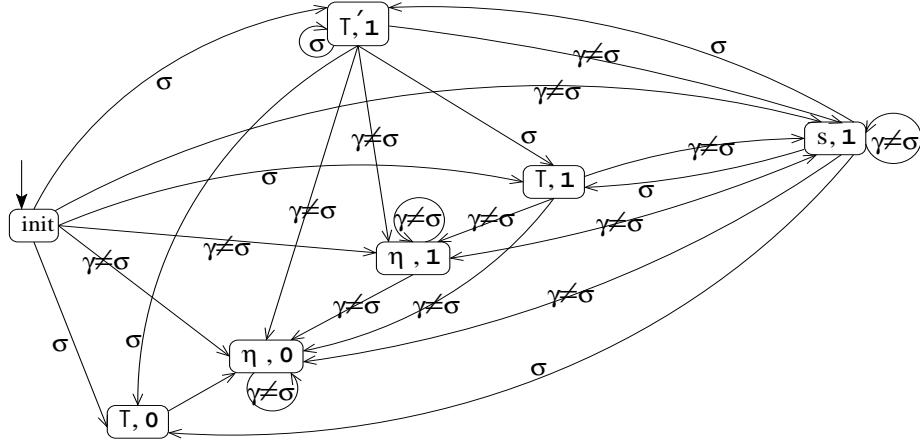
- Automaton for the query $/**[parent::\sigma]$



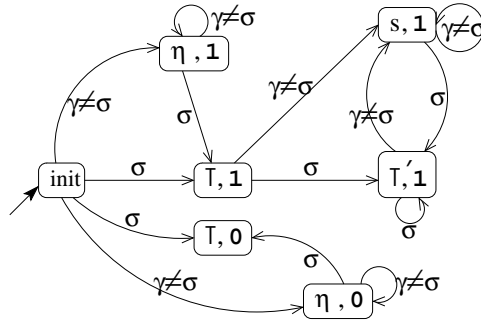
- Automaton for the query $/**[ancestor::\sigma]$



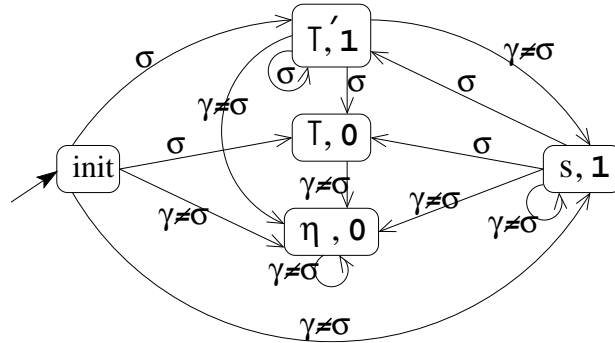
- Automaton for the query $//*[child::\sigma]$



- Automaton for the query $//*[preceding-sibling::\sigma]$



- Automaton for the query $//*[descendant::\sigma]$



A few words on some of the automata by way of explanation. First, the reason why the automaton for **self** does *not* have the states $(T, 0)$, $(T, 1)$, $(s, 1)$: for $(T, 0)$, $(T, 1)$, by the semantics of subsection 4.1 we must have $x = \sigma$, where x is the name of the current node on t , but then the query $//*[self::\sigma]$ should select the current node, so one cannot be at such a state; as for $(s, 1)$, the reasoning is just the opposite. Next, the reason why the automaton for **descendant** does *not* have the states $(\eta, 1)$, $(T, 1)$: if the semantics attribute one of these pairs to any node u , that would mean the node u has a selected descendant u' ; which means that u' has some σ -descendant node, which would then be a σ -descendant for u too, so Q should select u .

5 Maximal Priority Runs of Basic Query Automata

Note that the following properties, required by our semantics of subsection 4.1, hold on the automata \mathbf{A}_Q constructed above, for any basic Core XPath query $Q = //*[axis::\sigma]$:

- i) There are no transitions from any state with boolean component 0 to a state with boolean component 1;
- ii) The σ -transitions have all their target states in $\{(\top, 0), (\top, 1), (\top', 1)\}$; and for any $\gamma \neq \sigma$, the target states of γ -transitions are all in $\{(\eta, 0), (\eta, 1), (s, 1)\}$.

Theorem 1 *Let Q be any basic Core XPath query, t any given trdag, and let \mathcal{G} denote either the rlag \mathcal{D}_t , or any given chibling \mathcal{F} of \mathcal{L}_t . Assume given a labeling function L from $\text{Nodes}(\mathcal{G})$ into the set of ll-pairs, which is correct with respect to Q , i.e., in conformity with the semantics of subsection 4.1. Then there is a run r of the automaton A_Q on \mathcal{G} , such that :*

- i) r is compatible with L ; i.e., $r(u) = L(u)$ for every node u on \mathcal{G} ;
- ii) r satisfies the maximal priority condition (**MP**) of subsection 4.1.

Proof. We first construct, by induction, a ‘complete’ run (i.e., defined at all the nodes of \mathcal{G}) satisfying property i). For that, we shall employ reasonings that will be specific to the axis of the basic query Q . We give here the details only for the axis **parent**; they are similar for the other axes.

$Q = //*[parent::\sigma]$: (The axis considered is non-sibling so $\mathcal{G} = \mathcal{D}_t$ here.) At the root u node of \mathcal{D}_t , we set $r(u) = L(u)$; we have to show that there is a transition rule in \mathbf{A}_Q of the form $(init, llab(u)) \rightarrow L(u)$. Obviously, for the axis **parent**, the root node u cannot correspond to a node on t selected by Q , so the only ll-pairs possible for $L(u)$ are $(l, 0), (l, 1)$, with $l \in \{\eta, \top\}$; for each of these choices, we do have a transition rule of the needed form, on \mathbf{A}_Q .

Consider then a node v on \mathcal{D}_t such that, at each of its ancestor nodes u on \mathcal{D}_t , the part of the run r of A_Q has been constructed such that $r(u) = L(u)$; assume that the run cannot be extended at the node by setting $r(v) = L(v)$. This means that there exists a parent node w of v , such that $(L(w), llab(v)) \rightarrow L(v)$ is not a transition rule of \mathbf{A}_Q ; we shall then derive a contradiction. We only have to consider the cases where the boolean component of $L(w)$ is greater than or equal to that of $L(v)$. The possible couples $L(w), L(v)$ are then respectively:

$$\begin{aligned} L(w) &: (\top, 0) \mid (\top, 1) \mid (\top, 1) \mid (\top', 1) \mid (\top', 1) \\ L(v) &: (\eta, 0) \mid (\top, 1) \mid (\eta, 1) \mid (\top, 1) \mid (\eta, 1) \end{aligned}$$

In all cases, we have $llab(w) = \sigma$ because of the semantics, so the node (on t corresponding to the node) v has a σ -parent, so must be selected; thus the above choices for $L(v)$ are not in conformity with the semantics; contradiction.

We now prove that the complete run r thus constructed, satisfies property ii). For this part of the proof, the reasoning does not need to be specific for each Q ; so, write Q more generally, as $//*[axis::\sigma]$ for some given σ . Suppose the run r does not satisfy the maximal priority condition at some node v on \mathcal{G} ; assume, for instance, that the run r made the choice, say of the ll-pair $(l, 1)$, although the maximal labeling of the node v , in a manner compatible with the ll-pairs of all its parents, was the ll-pair $(l, 0)$. Since L is assumed correct, and r is compatible with L , the maximal possible labeling $(l, 0)$ would mean that the node (on t corresponding to the node) v has no descendant selected by Q ; whereas, the choice that r is assumed to have made at v , namely the ll-pair $(l, 1)$, has the opposite semantics whether or not $llab(v) = \sigma$; in other words, the labeling L would not be correct with respect to Q ; contradiction. The other possibilities for the ‘bad’ labelings under r also get eliminated in a similar manner. \square

Theorem 2 *Let $Q, t, \mathcal{D}_t, \mathcal{F}, \mathcal{G}$ be as above. Let r be a (complete) run of the automaton A_Q on \mathcal{G} , which satisfies the maximal priority condition (**MP**) of*

subsection 4.1. Then the labeling function L on $\text{Nodes}(\mathcal{G})$, defined as $L(u) = r(u)$ for any node u , is correct with respect to the semantics of subsection 4.1.

Proof. Let us suppose that the labeling L deduced from r is *not* correct with respect to Q ; we shall then derive a contradiction. The reasoning will be by case analysis, which will be specific to the axis of the basic query Q considered. We give the details here for $Q = // * [\text{descendant} : : \sigma]$. The axis is non-sibling, so we have $\mathcal{G} = \mathcal{D}_t$ here. The sets $\text{Nodes}(t), \text{Nodes}(\mathcal{D}_t)$ are in a natural bijection, so for any node u on \mathcal{D}_t we shall also denote by u the corresponding node on t , in our reasonings below.

We saw that the automaton \mathbf{A}_Q for the **descendant** axis does not have the states $(\eta, 1), (\top, 1)$. Consider then a node u on \mathcal{D}_t such that: for all ancestor nodes w of u , the ll-label $r(w)$ is in conformity with the semantics, but the ll-pair $r(u)$ is not in conformity. Now, \mathbf{A}_Q has only 5 states: $(\text{init}), (\top', 1), (s, 1), (\top, 0), (\eta, 0)$, of which only the last four can ll-label the nodes. So the possible ‘bad’ choices that r is assumed to have made at our node u , are as follows:

- (a) $r(u) = (\top', 1)$, but the node u is *not* an answer to the query Q . Here $\text{name}(u)$ must be σ , so the choice of r ought to have been $(\top, 0)$;
- (b) $r(u) = (s, 1)$, but the node u is *not* an answer to the query Q . Here $\text{name}(u) \neq \sigma$, so the choice of r ought to have been $(\eta, 0)$;
- (c) $r(u) = (\eta, 0)$, but the node u *is* an answer to the query Q . Here $\text{name}(u) \neq \sigma$, so the choice of r ought to have been $(s, 1)$;
- (d) $r(u) = (\top, 0)$, but the node u *is* an answer to the query Q . Here $\text{name}(u)$ must be σ , so the choice of r ought to have been $(\top', 1)$.

In all the four cases, we have to show:

- i) that the “ought-to-have-been” choice ll-pair is reachable from *all* the parent nodes of u ;
- ii) *and* that, with such a new and ‘correct’ choice made at u , r can be completed from u , into a run on the entire dag \mathcal{D}_t .

The reasoning will be similar for cases (a), (b), and for the cases (c), (d). Here are the details for case (a): That u is *not* an answer to Q means that u has *no* σ -descendant node, so for all nodes v below u on \mathcal{D}_t , we have $\text{llab}(v) \neq \sigma$. Therefore, assertions i) and ii) above follow from the following observations on the automaton for $Q = // * [\text{descendant} : : \sigma]$:

- i) *if* r could reach the state $(\top', 1)$ at node u (via a σ -transition) from any parent node of u , then $(\top, 0)$ is also reachable thus at u , from any of them;
- ii) *if*, from the state $(\top', 1)$, r could reach all the nodes on \mathcal{D}_t below u (with state $(\eta, 0)$), via transitions over $\gamma \neq \sigma$, then it can do exactly the same now, with the ‘correct’ choice ll-pair $(\top, 0)$ at u .

As for case (c): Node u *is* an answer to Q here, so u has a σ -descendant; let v be a σ -node below u on \mathcal{D}_t ; the ll-pair $r(v)$ that r assigns to v must then be either $(\top', 1)$ or $(\top, 0)$; this implies that r passed from the state $(\eta, 0)$ – supposedly assigned by r to u – to $(\top', 1)$ or $(\top, 0)$ somewhere between u and v ; which is impossible, as is easily seen on the automaton \mathbf{A}_Q for the axis **descendant** considered. The reasoning for case (d) is even easier: from state $(\top, 0)$, no state with an outgoing σ -transition is reachable. \square

6 Evaluating Composite Queries

A composite query is a query in standard form, but is not basic. We propose to evaluate such a query incrementally. For this, it suffices to consider queries that are of the form $// * [\mathbf{A} : : x \text{ conn } \mathbf{A}' : : x']$, where $\text{conn} \in \{\text{and}, \text{or}\}$, or of the form $// * [\mathbf{A}_1 : : * [\mathbf{A}_2 : : \sigma]]$. For those of the former type, we observe first that the components in a disjunction (resp. conjunction) under a ‘*’ can be evaluated separately. Indeed, the answer for $Q = // * [\mathbf{A} : : x \text{ conn } \mathbf{A}' : : x']$ can

be obtained as union (resp. intersection) of the answers for the two “component” queries $//*[A::x]$, and $//*[A'::x']$, when *conn* is an *or* (resp. an *and*). We apply the method described earlier, separately for $Q_1 = //*[A::x]$ and for $Q_2 = //*[A'::x']$, thus getting two respective evaluating runs r_1, r_2 . Any node u of the dag \mathcal{D}_t will then be re-labeled, by the composite query Q , with ll-pairs computed by a function **AND** when *conn* = *and* (resp. **OR** when *conn* = *or*), in conformity with the semantics presented in the Section 4.1:

$$\begin{aligned} AND(u) &= (s, 1) \text{ iff } r_1(u) = (l', 1) = r_2(u); \\ AND(u) &= (\eta, 0) \text{ iff } r_1(u) = (l, 0) \text{ or } r_2(u) = (l, 0); \\ AND(u) &= (\eta, 1) \text{ otherwise.} \\ OR(u) &= (s, 1) \text{ iff } r_1(u) = (l', 1) \text{ or } r_2(u) = (l', 1); \\ OR(u) &= (\eta, 0) \text{ iff } r_1(u) = (l, 0) = r_2(u); \\ OR(u) &= (\eta, 1) \text{ otherwise.} \end{aligned}$$

Figure 5 below illustrates the above reasoning, for the evaluation of the composite query $Q = //*[\text{self}::b \text{ and } \text{parent}::a]$, on the trdag t of Figure 3:

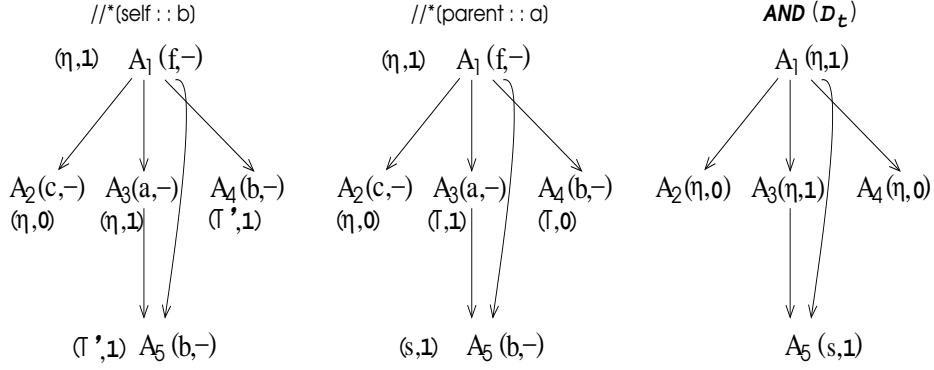


Fig. 5.

We next consider the queries of the form $Q = //*[A_1::*[A_2::\sigma]]$, with imbricated predicates. For their evaluation, we first consider a maximal priority run evaluating r_2 (resp. a set of runs \hat{r}_2) of the automaton associated to the inner query $//*[A_2::\sigma]$, on \mathcal{D}_t (resp. the set of all chiblings of \mathcal{L}_t). This run (resp. set of runs) will output the rag $r_2(\mathcal{D}_t)$ (resp. $\hat{r}_2(\mathcal{D}_t)$), as described in Section 4.2. Evaluating the imbricated query Q on the dag t is then done by running the automaton for the basic outer query $//*[A_1::s]$ on $r_2(\mathcal{D}_t)$ (resp. $\hat{r}_2(\mathcal{D}_t)$).

Finally, the answer for a query of the type $Q = //*[\text{child}::*[position() = k]]$, is the subset of the nodes answering $//*[\text{child}::*]$, which correspond to a k -th node on some chibling.

7 Conclusion

Our concern in this paper has been two-fold. The first part addressed the problem of running any bottom-up (unranked) tree automaton indifferently on a tree or on any of the dags obtained from the tree by full or partial compression; this gave rise to the notions of Tree/Dags (trdags) and of Tree/Dag automata. The second part of the paper addressed the issue of retrieving information from a trdag representing an XML document possibly given in a compressed form. (Note: Information retrieval from compressed structures, without having to uncompress them, is a field of active research; cf. e.g., [17, 11].) Limiting our concern here to the evaluation of queries formulated in terms of XPath axes, and more precisely to positive Core XPath queries, we have presented a method for evaluating them on any trdag t , *without* having to uncompress t , by breaking up the given query

into sub-queries of a basic type; with each basic query, an automaton is associated such that an unambiguous maximal priority run of this automaton can evaluate the query. An algorithm constructing these maximal priority runs is given in Appendix II; it has just been implemented (in Java). It is of complexity $\mathcal{O}(n^3)$ where n is the number of nodes of the given trdag t ; the bound $\mathcal{O}(n^3)$ is due to the relation *Parents* –less trivial on a trdag than on a tree; the complexity reduces to $\mathcal{O}(n^2)$ on trees.

One apparent advantage of the approach presented here is that the various basic sub-queries “composing” a given query can be evaluated in parallel, in several cases; a detailed analysis of this issue could be a possible direction for future work. A second possible direction would be to see if the Core XPath query evaluation algorithms of [10] can be adapted to dags, in a manner compatible with our approach. We also expect to be able to extend our approach to the evaluation of more general XPath queries, such as those involving the data values, by suitably adapting its underlying mechanism based on labeling.

References

1. S. Anantharaman, P. Narendran, M. Rusinowitch, *Closure Properties and Decision Problems of Dag Automata*, In Information Processing Letters, 94(5):231–240, 2005.
2. P. Buneman, M. Grohe, C. Koch, *Path queries on compressed XML*. In Proc. of the 29th Conf. on VLDB, 2003, pp. 141–152, Ed. Morgan Kaufmann.
3. G. Busatto, M. Lohrey, S. Maneth, *Grammar-Based Tree Compression*. EPFL Technical Report IC/2004/80, <http://icwww.epfl.ch/publications>.
4. G. Busatto, M. Lohrey, S. Maneth, *Efficient Memory Representation of XML Documents*. In Proc. DBPL’05 (to appear), LNCS 3774, Springer-Verlag, 2005.
5. W. Charatonik, *Automata on DAG Representations of Finite Trees*, Technical Report MPI-I-99-2-001, Max-Planck-Institut für Informatik, Saarbrücken, Germany.
6. H. Comon, M. Dauchet, R. Gilleron, F. Jacquemard, D. Lugiez, S. Tison, M. Tommasi, *Tree Automata Techniques and Applications*, <http://www.grappa.univ-lille3.fr/tata/>
7. M. Frick, M. Grohe, C. Koch, *Query Evaluation of Compressed Trees*, In Proc. of LICS’03, IEEE, pp. 188–197.
8. G. Gottlob, C. Koch, *Monadic Queries over Tree-Structured Data*, In Proc. of LICS’02, IEEE,
9. G. Gottlob, C. Koch, *Monadic Datalog and the Expressive Power of Languages for Web Information Extraction*, In Journal of the ACM, 51(1):12–28, 2004.
10. G. Gottlob, C. Koch, R. Pichler, L. Segoufin, *The complexity of XPath query evaluation and XML typing* In Journal of the ACM 52(2):284–335, 2005.
11. M. Lohrey, *Word problems and membership problems on compressed words* In SIAM Journal of Computing, 35(5):1210–1240, 2006.
12. M. Marx, *XPath and Modal Logics for Finite DAGs*. In Proc. of TABLEAUX’03, pp. 150–164, LNAI 2796, 2003.
13. W. Martens, F. Neven, *On the complexity of typechecking top-down XML transformations*, In Theoretical Computer Sc., 336(1): 153–180, 2005.
14. M. Murata, A. Tozawa, M. Kudo, *XML Access Control Using Static Analysis*, In Proc. of the 10th ACM Conf. on Computer and Communications Security (CCS’03), pp.73–84, ACM, 2003.
15. F. Neven, *Automata Theory for XML Researchers*, In SIGMOD Record 31(3), September 2002.
16. F. Neven, T. Schwentick, *Query automata over finite trees*, In Theoretical Computer Science, 275(1–2):633–674, 2002.
17. W. Rytter, *Compressed and fully compressed pattern matching in one and two dimensions*, In Proceedings of the IEEE, 88(11):1769–1778, 2000.
18. J.W. Thatcher, J.B. Wright, *Generalized finite automata theory with an application to a decision problem of second-order logic*, In Math. Syst. Theory, 2(1):57–81, 1968.
19. Worl Wide Web Consortium, *XML Path Language (XPath Recommendation)*, <http://www.w3c.org/TR/xpath/>

Appendix I: From Canonical Forms to Standard Forms

We stick to the notations of Section 3.1. Given any canonical XPath expression Q_{can} , we compute, inductively, an equivalent standard XPath expression denoted as $Std(Q_{can})$; as earlier, *conn* stands for either of the boolean connectives *and*, *or*.

To start with, we define:

```

Std([true]) = self::*
Std([Scan]) = Scan
Std([A:: $\sigma$ [Scan]]) = A::*[self:: $\sigma$  and Std([Scan])]
Std([A:: $\sigma$ [A1:: $\sigma_1$ [... [Ak:: $\sigma_k$ ]...]]) = A::*[self:: $\sigma$  and
    A1::*[self:: $\sigma_1$  and... Ak-1::*[self:: $\sigma_{k-1}$  and Ak:: $\sigma_k$ ] ...]]

```

We also define, for every basic axis relation, an inverse relation, as follows:

```

self-1 = self
child-1 = parent
parent-1 = child
ancestor-1 = descendant
descendant-1 = ancestor
following-sibling-1 = preceding-sibling
preceding-sibling-1 = following-sibling

```

For any query $Q = // * [X]$ in standard form, we set $exp(Q) = X$. For any canonical XPath query $Q = /C_1/C_2/\dots/C_n$, the standard form $Std(Q)$ of Q is then generated by the following recursive construction:

Case of length 1: $Q = /C_1$

```

Std(/child:: $\sigma$ [Xcan]) = //*[(Root and self:: $\sigma$ ) and Std([Xcan])]
Std(/child::*[Xcan]) = //*[Root and Std([Xcan])]
Std(/descendant:: $\sigma$ [Xcan]) = //*[self:: $\sigma$  and Std([Xcan])]
Std(/descendant::*[Xcan]) = //*[Std([Xcan])]
Std(/axis:: $\sigma$ [Xcan] conn axis':: $\sigma'$ [X'can]) =
    //*[exp(Std(/axis:: $\sigma$ [Xcan])) conn exp(Std(/axis':: $\sigma'$ [X'can]))]

```

Case of length n>1: $Q = /C_1/C_2/\dots/C_n$

```

Std(/C1/.../Cn-1/A:: $\sigma$ [Xcan]) =
    //*[(self:: $\sigma$  and Std([Xcan])) and A-1::*[exp(Std(/C1/.../Cn-1))]
Std(/C1/.../Cn-1/A:: $\sigma$ [Xcan] conn A':: $\sigma'$ [X'can]) =
    //*[((self:: $\sigma$  and Std([Xcan])) conn (self:: $\sigma'$  and Std([X'can]))
        and A-1::*[exp(Std(/C1/.../Cn-1))]

```

This translation procedure is of complexity linear with respect to the total number of location steps (i.e. of the form $axis:: σ) that appear in Q .$

Appendix II: Constructing the Maximal Priority Run

Given a trdag t , and a basic non-sibling query Q we give here the algorithm constructing the maximal priority run of the automaton \mathbf{A}_Q on \mathcal{D}_t . (It is trivial for the sibling queries.) Let n be the number of nodes on \mathcal{D}_t ; The idea is to construct a directed acyclic graph $\mathbb{G} = (\mathbb{V}, \mathbb{E})$, representing all possible (complete) runs of \mathbf{A}_Q on \mathcal{D}_t , and at the same time choose the one of maximal priority. The set \mathbb{V} of vertices of \mathbb{G} are elements of the form $A_{i,(t,x)}$, $1 \leq i \leq n$, where A_i is a non-terminal of \mathcal{L}_t , and (t,x) is an ll-pair such that (t,x) is in $States(\mathbf{A}_Q)$ such that there exists a (complete) maximal priority run r of \mathbf{A}_Q on \mathcal{D}_t with $r(A_i) = (t,x)$. The algorithm has four steps as described below (where $llab(A_i)$ is the first component of the current label at the node on \mathcal{D}_t named A_i).

Step i): For every non-terminal A_i of \mathcal{L}_t , $1 \leq i \leq n$ compute the set $Parents(A_i)$ of its parent non-terminals. To render the presentation uniform, we shall introduce a ‘fictive’ symbol A_0 , and set $Parents(A_1) = \{A_0\}$.

The cost of this step is $\mathcal{O}(n^3)$: For every A_i we have to check if A_i is in $Sons(A_j)$ for $j \in \{1, \dots, i-1\}$. The maximal size of $Sons(A_j)$ can be $(n-j)$, thus, for A_i we have $\sum_{j=1}^{i-1} (n-j) = (i-1)n - \sum_{j=1}^{i-1} j \in \mathcal{O}(i(n-i))$. In total we get: $\sum_{i=1}^n i(n-i) \leq n \sum_{i=1}^n i \in \mathcal{O}(n^3)$.

Step ii): Construct the set of vertices and the arcs of the graph $\mathbb{G} = (\mathbb{V}, \mathbb{E})$.

For every $i \in \{1, \dots, n\}$ we have at most three different (t,x) determined by $llab(A_i)$, so constructing the vertices costs $\mathcal{O}(3n) = \mathcal{O}(n)$.

Next, a pair $(A_{j,(t',x')}, A_{i,(t,x)})$ will be an arc in \mathbb{E} iff $A_{j,(t',x')}, A_{i,(t,x)}$ are in \mathbb{V} , and the rule $((t',x'), llab(A_i)) \rightarrow (t,x)$ is a transition of \mathbf{A}_Q . So, for every given vertex $A_{j,(t',x')} \in \mathbb{V}$, we construct the set $Target(A_{j,(t',x')})$ of indices of all non-terminals A_i , reachable (with some label), from $A_{j,(t',x')}$ on \mathbb{G} . The cost of this construction is in $\mathcal{O}(n^2)$: indeed, for every vertex $A_{j,(t',x')}$ and for every $A_{i,(t,x)} \in \mathbb{V}$ such that $A_j \in Parents(A_i)$, we have to check if the former is reachable from the latter; that gives us $3\#(Parents(A_i))$ transitions of \mathbf{A}_Q to check, so, on the whole, we get: $3 \sum_{i=1}^n 3\#(Parents(A_i)) = 9 \sum_{i=1}^n (i-1) \in \mathcal{O}(n^2)$.

Step iii): Eliminate the incomplete runs.

Consider any two given $i, j \in \{1, \dots, n\}$, and suppose $A_{j,(t',x')}$ is a vertex such that $j \in Parents(A_i)$ but $i \notin Target(A_{j,(t',x')})$; this means A_i is not reachable from the former vertex if we associate the label (t',x') to A_j ; and this implies that the vertex $A_{j,(t',x')}$ is incorrect; we remove from \mathbb{G} all incorrect vertices, and all the incoming arcs at such vertices (by a bottom-up traversal of \mathbb{G}).

Finding all the incoming arcs at a given $A_{i,(t,x)}$, costs $3\#Parents(A_i) = 3(i-1)$ checks; so, the total cost of this step amounts to $3 \sum_{i=1}^n 3(i-1) \in \mathcal{O}(n^2)$.

Step iv): Choose the maximal priority run.

Begin by setting $r(A_0) = init$; then, for every A_i , $i \in \{1, \dots, n\}$, associate a pair (t,x) such that: $A_{i,(t,x)} \in \mathbb{V}$ and (t,x) is of maximal priority between all possible pairs (t',x') such that $A_{i,(t',x')} \in \mathbb{V}$. The cost of this last step is in $\sum_{i=1}^n 3(i-1) \in \mathcal{O}(n^2)$.

The total cost of the algorithm is therefore $\mathcal{O}(n^3)$.

Given:

The dependency dag \mathcal{D}_t , $n =$ its number of nodes;
 $A_i, 1 \leq i \leq n$: the non-terminals of \mathcal{L}_t naming the nodes of \mathcal{D}_t ;
 \mathbf{A}_Q the automaton for the basic query Q ,
 $States =$ the set of states of \mathbf{A}_Q , $\Delta =$ the set of transition rules of \mathbf{A}_Q .

BEGIN :

```

/* Step i): Construct the relation Parents */
Parents( $A_1$ ) :=  $\{A_0\}$ ; for all  $i \in \{2, \dots, n\}$  Parents( $A_i$ ) :=  $\emptyset$ ;
For  $i \in \{2, \dots, n\}$  do
  For  $j \in \{1, \dots, i-1\}$  do
    If  $A_i \in Sons(A_j)$  Parents( $A_i$ ) := Parents( $A_i$ )  $\cup \{A_j\}$ ;
  od; od;

/* Step ii): Construct the vertices  $\mathbb{V}$ , the arcs  $\mathbb{E}$ , and the Target sets */
 $\mathbb{V} := \{A_{0,init}\} \cup \{A_{i,(t,x)} \mid 1 \leq i \leq n, t \text{ defined by } llab(A_i), (t,x) \in States\}$ 
 $\mathbb{E} := \emptyset$ ; For all  $A_{i,(t,x)} \in \mathbb{V}$  Target( $A_{i,(t,x)}$ ) :=  $\emptyset$ ;
For all  $A_{i,(t,x)} \in \mathbb{V}$ , with  $i \geq 1$ 
  For all  $A_j \in Parents(A_i)$ 
    If  $A_{j,(t',x')} \in \mathbb{V}$  and  $((t',x'), llab(A_i)) \rightarrow (t,x) \in \Delta$  {
       $\mathbb{E} := \mathbb{E} \cup \{(A_{j,(t',x')}, A_{i,(t,x)})\}$ ;
      Target( $A_{j,(t',x')}$ ) := Target( $A_{j,(t',x')}$ )  $\cup \{i\}$ ;
    }
  }

/* Step iii): Eliminate the incomplete runs, via a bottom-up search */
i :=  $n$ ;
(#) For all  $A_{i,(t,x)}, A_{i,(t',x')} \in \mathbb{V}$ 
  If Target( $A_{i,(t,x)}$ )  $\subsetneq$  Target( $A_{i,(t',x')}$ ) {
     $\mathbb{V} := \mathbb{V} \setminus \{A_{i,(t,x)}\}$ ;
    For all  $A_{k,(t',x')} \in \mathbb{V}$  such that  $(A_{k,(t',x')}, A_{i,(t,x)}) \in \mathbb{E}$  {
       $\mathbb{E} := \mathbb{E} \setminus \{(A_{k,(t',x')}, A_{i,(t,x)})\}$ ;
      If  $\mathbb{E} \cap \{(A_{k,(t',x')}, A_{i,(t,y)}) \mid (t,y) \text{ defined by } llab(A_i)\} = \emptyset$ 
        Target( $A_{k,(t',x')}$ ) := Target( $A_{k,(t',x')}$ )  $\setminus \{i\}$ ;
      }
    }
  }
i :=  $i-1$ ;
If  $i \geq 1$  GOTO (#);

/* Step iv): Construct now the maximal priority run, top-down */
 $r(A_0) := init$ ;  $i := 1$ ;
($)  

 $r(A_i) := \max\{(t,x) \mid A_{i,(t,x)} \in \mathbb{V}\}$ ;  

If  $i < n$ ,  $\{i := i+1$ ; GOTO ($) $\}$ ;  

Else RETURN  $\{r(A_1), \dots, r(A_n)\}$ ;  

END.

```

Appendix III: A Complete Example

1) We evaluate the query $Q = /descendant::*[descendant::b [parent::a]]$ on the partially compressed document t , given to the left of Figure 6. Note that we want to select every node having some descendant b with parent a . To start with, we first translate Q into standard form, as:

$$Q = //*[descendant::*[self::b \text{ and } parent::a]]$$

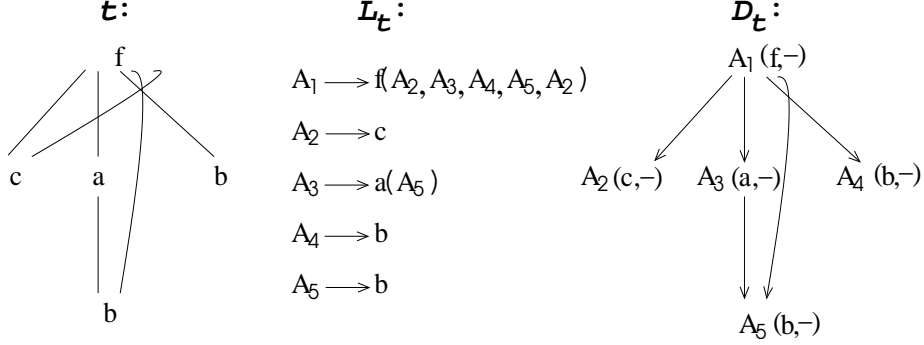


Fig. 6. Document t , its normalized Grammar \mathcal{L}_t , and the Dependency rlag \mathcal{D}_t

Figure 7 represents, to the left, the rlag \mathcal{D}_t labeled by the automaton for the query $//*[self::b]$; to the middle, the same rlag labeled by the automaton for the query $#[parent::a]$; and to the right, the rlag \mathcal{D}'_t re-labeled for the conjunction, as explained in Section 6.

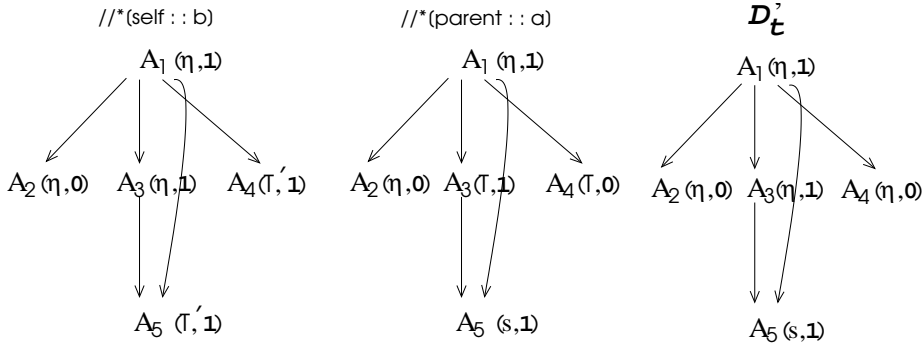


Fig. 7.

Figure 8 shows, to the left, the rlag obtained by re-labeling \mathcal{D}'_t with the run of the automaton for the query $#[descendant::s]$; and to the right, the (minimal) sub-rlag of \mathcal{D}_t formed of nodes marked now by ll-pairs with boolean component 1, and the corresponding answer for the query Q on the document t . This final sub-rlag of \mathcal{D}_t is obtained by cutting out the nodes where the ll-pairs attached have boolean component 0.

2) On the same document t as above, we consider now the following standard form query: $Q' = //*[child::b \text{ or } following-sibling::b]$

To the left of Figure 9 is the rlag \mathcal{D}_t labeled by the run of the automaton for $#[child::b]$; and to the right, the labeling of the 3 chiblings of \mathcal{D}_t , by the run of the automaton for $#[following-sibling::b]$.

The two rlags $\mathcal{D}'_t, \mathcal{D}''_t$ of Figure 10 below, are then obtained by applying the re-labeling functions respectively lab_r and λ_r (of subsection 4.2) for these

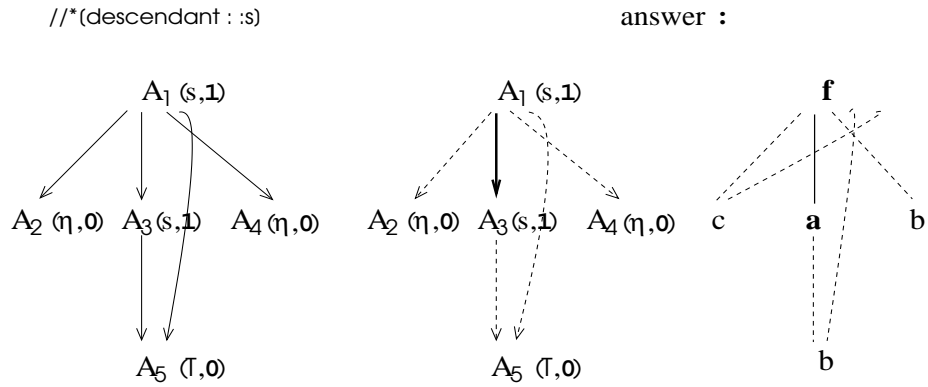


Fig. 8.

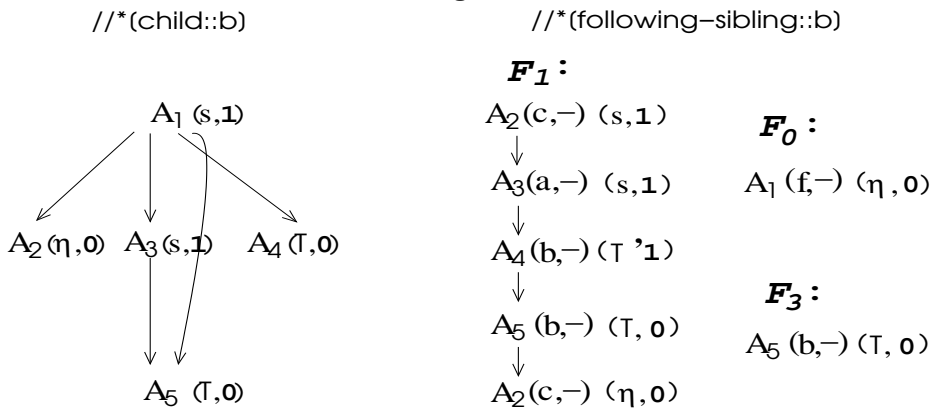


Fig. 9.

respective runs. The rag D_t''' to the right is obtained by the run of automaton for `/** [self::s]` on D_t'' .

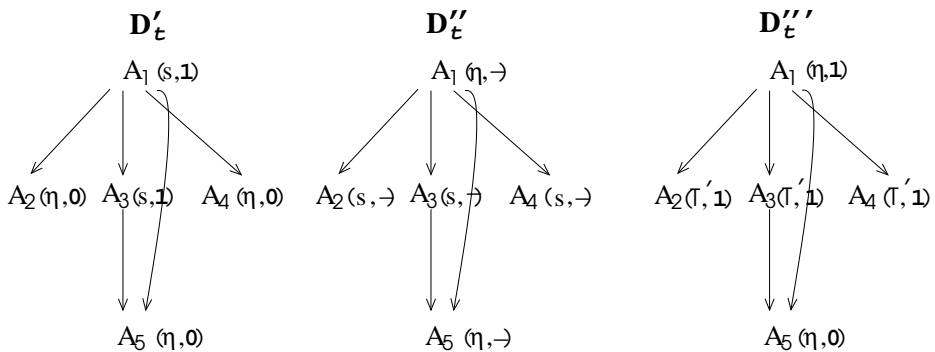


Fig. 10.

The Figure 11 presents the final answer for our query obtained by applying the function OR from section 6.

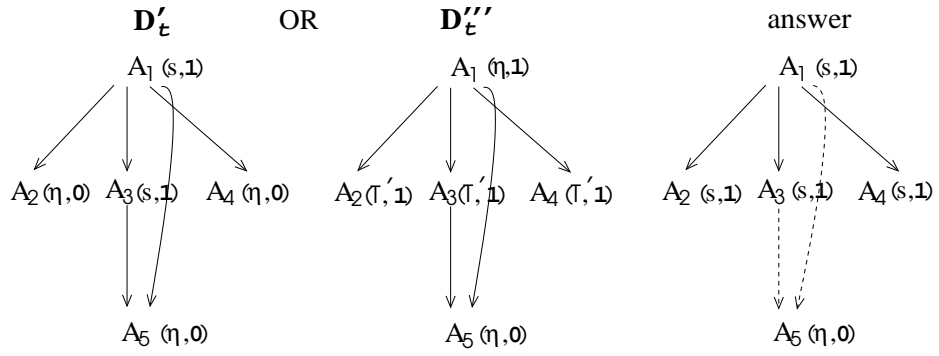


Fig. 11.