

# View-Based Tree-Language Rewritings for XML

Laks V.S. Lakshmanan<sup>1</sup> and Alex Thomo<sup>2</sup>

<sup>1</sup> University of British Columbia, Vancouver, BC, Canada, [laks@cs.ubc.ca](mailto:laks@cs.ubc.ca)

<sup>2</sup> University of Victoria, Victoria, BC, Canada, [thomo@cs.uvic.ca](mailto:thomo@cs.uvic.ca)

**Abstract.** We study query rewriting using views (QRV) for XML. Our queries and views are regular tree languages (RTLs) represented by tree automata over marked alphabets, where the markers serve as “node selectors”. We formally define query rewriting using views for RTLs and give an automata-based algorithm to compute the maximally contained rewriting. The formalism we use is equal in power with Monadic Second Order (MSO) logic, and our algorithm for computing QRV is the first to target this expressive class. Furthermore we prove a tight lower bound, thus showing that our algorithm is optimal. Another strength of our automata-based approach is that we are able to cast computing QRV into executing a sequence of intuitive operations on automata, thus rendering our approach practical as it can be easily implemented utilizing off-the-shelf automata toolboxes. Finally, we generalize our framework to account for more complex queries in the spirit of the **FOR** clause in XQuery. For this generalization as well, we give an optimal algorithm for computing the maximally contained rewriting of queries using views.

## 1 Introduction

Query rewriting using views (QRV) is a fundamental problem that finds wide applications in query optimization, data integration, data warehousing, security, and other critical database services. In this paper we study QRV for node-selecting queries and views, over XML trees. As there are several variants of views for answering queries we begin by illustrating the classical QRV problem we focus on in this paper.

*Example 1.* Suppose we have a very large collection of movies such as *imdb.com*, organized in a super-tree, with each movie being a sub-tree containing title, year, and characters. Suppose the character nodes branch out into actors playing the character. Sometimes, a movie contains characters played by more than one actor. We call these characters “multi-actor characters (MAC)”. Consider a query which returns all the movie sub-trees having a MAC. Such movies are a small minority; their number is about 50. Assume that the result of this query is materialized into a view. Obviously this view is of tremendous help in answering some new queries, such as “find all the actors playing a MAC”. We can rewrite this new query into “find all the actors playing a MAC in a movie having a MAC” and answer it on the view-extension instead of accessing the original database. The difference in performance is huge; using the view materialization takes 50

accesses, whereas using the original database takes hundreds of thousands of accesses (there are about 700,000 movies as per [imdb.com](http://imdb.com)).  $\square$

Due to its importance, QRV has been explored for fragments of XPath (cf. [5, 31, 16, 8, 2, 30]) and fragments of XQuery (cf. [32, 24, 4]). In this paper, we study the problem for queries and views represented by tree automata, which provide for a more general approach that can elegantly capture fine grained structure along vertical and horizontal axes. These queries can be easily specified using a DTD-like syntax, rendering them user-friendly devices for querying XML. While the techniques and tools introduced in the aforementioned works are interesting and well-founded for computing QRV for XPath and XQuery, they do not seem to extend to dealing with the more general setting of queries and views represented by tree automata, which have desirable properties with respect to expressivity in querying XML. In regard to expressivity, Neven and Schwentick [20, 21] and Schwentick [26] argue for formalisms as expressive as monadic second order logic (MSO) for specifying node-selecting queries. This is sometimes called a “golden standard” against which query formalisms should be measured up.

In this paper, we study the QRV problem for queries and views represented by automata equivalent in power to MSO. It is worth pointing out that being able to handle this target expressivity is one of the strengths of our approach to QRV. Another strength is the fact that computing QRV is cast into executing a sequence of intuitive operations on automata. This makes our solution quite practical, as it can be easily implemented on top of the readily available automata toolboxes (e.g. [12, 9]).

Our queries and views are first described as sets of tree-position pairs in order to facilitate the definitions and understanding of rewritings. In practice the queries and views are specified by finite tree automata over alphabets with marked symbols that serve the purpose of selecting nodes in XML trees. More specifically, for our constructions we use colors as markers, which makes the development easier. This formalism is equivalent to the formalism of querying trees using automata with boolean markings (ABM) [cf. [28], further developed in [23]].

Automata over colored alphabets provide us with critical advantages in computing QRV. The first advantage is the ability to express and construct a series of intermediate languages for obtaining query rewritings. Being able to use multiple colors for marking regions of interest in the trees of these intermediate languages is crucial to our approach. Our language constructions are also facilitated by the one-way nature of the automata we use. The second advantage is the ability to determinize the automata. This is key to our main construction for query rewriting.

Summarizing, we provide an algorithm for QRV for the general case where queries and views are implemented in an automata framework equal in power to MSO, to our knowledge, for the first time. Our algorithm runs in *singly-exponential time*. We show that QRV in our case is EXPTIME-hard, thus showing the optimality of our algorithm. Additionally, we show how to extend our results to reasoning about QRV for the more general case of queries returning

a forest of trees as an answer in the spirit of the `FOR` clause in XQuery. For this generalization as well, we provide an optimal algorithm to compute query rewritings using views.

We make the following contributions.

1. We define tree-pattern operators to cleanly express view-based rewriting as the solution to a view-query equation (Section 4).
2. Next, we present an algorithm for computing query rewritings using views. We define languages *over colored alphabets* and present a series of intermediate automata constructions, which we believe are of independent interest. We prove an EXPTIME lower bound for QRV, and show that our algorithm is optimal (Section 5).
3. We generalize our results to the case where queries select more than one subtree and produce a set of forests as output. We show that we are still able to compute view-based rewritings in singly-exponential time, thus being again optimal with respect to the above lower bound (Sections 6, 7, and 8).

## 2 Related Works

Automata theory has long been recognized as a useful tool for providing elegant solutions to challenging problems on XML (cf. [19, 21, 14, 22, 26, 17, 7]).

Two prominent automata-based approaches for querying XML data that have been proposed are (bottom-up) finite tree automata with selecting states (FTAS) introduced by Neven ([19], p. 128) and Frick et al. ([14]), and query automata (QA) introduced by Neven and Schwentick ([21]). Both formalisms are shown to capture precisely the queries definable in MSO and thus are the natural candidates for us to study the QRV problem in. However, two important issues prevent us from using either of these formalisms directly for our study of QRV: (1) Deterministic FTAS are too weak to express all queries expressible by non-deterministic FTAS; determinism is a key property that our techniques and algorithms rely on. (2) On the other hand, QA are deterministic, but two-way, and thus can go up and down a target tree multiple times, a feature that makes reasoning about QRV difficult.

Another nice standard for querying XML is Propositional Dynamic Logic (PDL) for trees ([1, 6]), which corresponds to Regular XPath. The latter was shown by ten Cate and Segoufin ([27]) to be not expressively complete for MSO. Regular XPath was extended later to  $\mu$ XPath by Calvanese et al. ([7]) to obtain the full power of MSO. It may be possible to use automata formalisms, such the one in [7] capturing  $\mu$ XPath, for QRV—this is an avenue we leave open for future exploration.

We close this section by briefly discussing two other works on query rewriting. Fan et al. ([11]) study rewritings of regular XPath queries defined over *virtual* views. This is a different problem from QRV that we consider here. The views in our case are materialized and the queries are over the original database. It would be interesting to see how our techniques could be used for the problem

studied by Fan et al., considering an alternate query formalism such as the one we propose that is complete for MSO.

Thomo and Venkatesh ([29]) use Visibly Pushdown Automata to model and rewrite (XML) schemas by using other schemas. Again, this problem is different from the problem we study in this paper.

### 3 Automata

We consider finite ordered trees – simply called *trees*. Also, we consider the trees to be *unranked*, which is to say that the nodes of the tree have an arbitrary (but finite) arity (cf. [10], p. 200). The nodes of the trees are labeled by symbols drawn from a fixed alphabet  $\Sigma$ . We denote by  $\mathcal{T}$  the set of all trees over  $\Sigma$ . We use  $a, b, \dots, e$  to denote labels in  $\Sigma$  and  $x, y, \dots$ , possibly with subscripts, to denote tree nodes. We denote by  $r_t$  the root of tree  $t$ , by  $\sigma_x$  the label of node  $x$ , and by  $N_t$  the set of nodes of a tree  $t \in \mathcal{T}$ .

**Definition 1.** A non-deterministic, bottom-up, finite tree automaton (FTA) over  $\Sigma$  is a quadruple  $\mathcal{A} = (S, \Sigma, F, \Delta)$ , where  $S$  is a finite set of states,  $F \subseteq S$  is a set of final states, and  $\Delta$  is a finite set of transition rules of the form  $H \xrightarrow{a} s$ , where  $H \subseteq S^*$  is a regular language over  $S$ ,  $a \in \Sigma$ , and  $s \in S$ .

Here,  $H$  is called a *horizontal* language. For simplicity, we blur the distinction between regular languages over  $S$  and the regular expressions used to specify them.

**Definition 2.** A tree  $t$  is accepted by an FTA  $\mathcal{A}$  if there exists a mapping  $\mu : N_t \rightarrow S$  such that:

1. If  $\mu(x) = s$ , then there is a transition rule  $H \xrightarrow{\sigma_x} s$  in  $\Delta$  with  $\mu(x_1) \dots \mu(x_n) \in H$ , where  $x_1, \dots, x_n$  are all the children of  $x$  in  $t$  in order.
2. If  $x$  is a leaf and  $\mu(x) = s$ , then there is a transition rule  $H \xrightarrow{\sigma_x} s$  in  $\Delta$  with  $\epsilon \in H$ .
3.  $\mu(r_t) \in F$ .

A mapping  $\mu$  as above specifies an *accepting run* of  $\mathcal{A}$  on  $t$ . An accepting run can be considered to be a tree of the same shape as  $t$  whose nodes are labeled by states given by the mapping  $\mu$ .

We denote by  $L(\mathcal{A})$  the set (language) of trees accepted by  $\mathcal{A}$ . A tree language  $L$  is said to be accepted (or recognized) by an FTA  $\mathcal{A}$  if  $L = L(\mathcal{A})$ . Tree languages recognized by FTAs are called *regular tree languages* (RTLs).

*Example 2.* Consider a collection of trees representing movies having, among other elements, one or more characters, each played by one or more actors. Let us use  $\mathbf{m}$ ,  $\mathbf{t}$ ,  $\mathbf{y}$ ,  $\mathbf{c}$ , and  $\mathbf{a}$  to abbreviate movie, title, year, character, and actor, respectively.

An FTA accepting movies having at least one character played by two or more actors is

$$\mathcal{A} = (\{s, s_{\mathbf{m}}, s_{\mathbf{c}}, s_{\mathbf{a}}\}, \Sigma, \{s_{\mathbf{m}}\}, \Delta)$$

where  $\Sigma \supset \{m, t, y, c, a\}$  and  $\Delta$  has the following transition rules:

$$\begin{aligned} s^* s_c s^* &\xrightarrow{m} s_m \\ s^* s_a s_a s^* &\xrightarrow{c} s_c \\ s^* &\xrightarrow{a} s_a \\ s^* &\xrightarrow{\Sigma} s. \end{aligned}$$

The last transition rule is a shorthand for saying that “we can go from  $s^*$  to  $s$  on any symbol of  $\Sigma$ .”

Let now  $t$  be the movie tree given in Figure 1, left. Clearly,  $t$  is accepted by  $\mathcal{A}$ . We assume that all the textual (unstructured data) have been identified (labeled) with a special symbol, say  $d \in \Sigma$ . An accepting run of  $\mathcal{A}$  on  $t$  is given in the same figure, right.

We describe in Section 4 how to modify  $\mathcal{A}$  to query for those movies having a character played by two or more actors.  $\square$

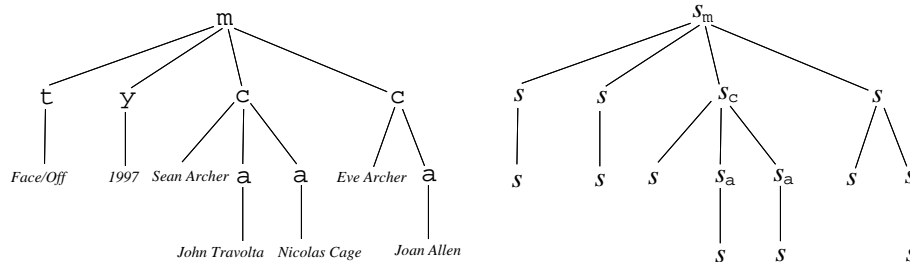


Fig. 1. A tree  $t$  and an accepting run of  $\mathcal{A}$  on  $t$ .

It is worth noting that an FTA can be specified using extended DTDs (EDTDs), a syntax that may be more familiar to a user (cf. [10], p. 233).

Consider again Definition 1. The FTA  $\mathcal{A}$  is called *deterministic* if  $H_1 \cap H_2 = \emptyset$  for all transitions  $H_1 \xrightarrow{a} s_1$  and  $H_2 \xrightarrow{a} s_2$ , where  $s_1 \neq s_2$ . In such a case, for any tree  $t$ , there can be at most one accepting run of  $\mathcal{A}$  on  $t$ .

FTAs that never get “stuck” are called *complete*. Given an FTA, it can be verified that the determinization procedure of [10] (p. 204) produces a deterministic FTA that is complete. Given a tree  $t$ , a complete and deterministic FTA (CDFTA)  $\mathcal{A}$  can read  $t$  in only one way, i.e., there is exactly one run of  $\mathcal{A}$  on  $t$ .

An FTA is *normalized* if for each symbol-state pair  $(a, s)$ , there is at most one transition  $H \xrightarrow{a} s$ . Each FTA can be transformed into a normalized FTA by unioning the left-hand sides of the transition rules labeled by the same symbol and having the same state on the right-hand side. If the FTA is deterministic, then it will remain so after this transformation.

*Intersection.* It is a well known fact that if  $L_1$  and  $L_2$  are RTLs, so is  $L_1 \cap L_2$ . However, the proof is typically done by encoding unranked trees into binary ones (cf. [10], p. 209). We can instead use a direct construction which preserves completeness and determinism. Given FTAs  $\mathcal{A}_1$  and  $\mathcal{A}_2$  for  $L_1$  and  $L_2$ , respectively, we can construct an FTA  $\mathcal{A}$  which recognizes  $L_1 \cap L_2$ . Furthermore, if  $\mathcal{A}_1$  and  $\mathcal{A}_2$  are complete and deterministic, so is  $\mathcal{A}$  (see Appendix).

### 3.1 Targetedness

Here we introduce our notion of “targeted” FTAs. Targeted FTAs turn out to be important in computing rewritings of queries using views.

**Definition 3.** *A normalized FTA is targeted if each state is the target (right side) of at most one transition rule.*

Any normalized FTA  $\mathcal{A} = (S, \Sigma, F, \Delta)$  can be transformed into a targeted one. For this, suppose that for a given state  $s$ , we have that

$$H_1 \xrightarrow{a_1} s, \dots, H_k \xrightarrow{a_k} s$$

are the transition rules in  $\Delta$  with target  $s$ . We introduce the operator  $\uparrow_s$  applied on  $\mathcal{A}$  with result  $\mathcal{A} \uparrow_s$  which is another FTA  $(S', \Sigma, F', \Delta')$  obtained from  $\mathcal{A}$  with the addition of fresh states  $s_1, \dots, s_k \notin S$ , and

$$\begin{aligned} S' &= (S \setminus \{s\}) \cup \{s_1, \dots, s_k\} \\ F' &= \begin{cases} (F \setminus \{s\}) \cup \{s_1, \dots, s_k\} & \text{if } s \in F \\ F & \text{otherwise} \end{cases} \\ \Delta' &= \{\gamma_s(H_i) \xrightarrow{a_i} s_i : i \in [1, k]\} \cup \\ &\quad \{\gamma_s(H) \xrightarrow{a} s' : H \xrightarrow{a} s' \text{ in } \Delta, \text{ and } s' \neq s\} \end{aligned}$$

where  $\gamma_s$  is substitution (in regular languages, c.f. [15], p. 60) of  $s$  by  $\{s_1, \dots, s_k\}$ . It can be verified that

**Lemma 1.**  $L(\mathcal{A}) = L(\mathcal{A} \uparrow_s)$ .

There is only a finite number of times—equal to  $|S|$ —we can apply  $\uparrow$  on a given FTA  $\mathcal{A}$  because the newly created states fulfill the targetedness condition. At the end of this procedure we obtain a targeted FTA, which we denote by  $\mathcal{A} \uparrow$ . From this, and Lemma 1, we have that

**Theorem 1.** *For each FTA  $\mathcal{A}$  we can obtain an equivalent targeted FTA  $\mathcal{A} \uparrow$  in PTIME.*

Furthermore, applying  $\uparrow$  on an FTA  $\mathcal{A}$ , preserves determinism, i.e. we have

**Proposition 1.** *If  $\mathcal{A}$  is deterministic, then so is  $\mathcal{A} \uparrow$ .*

## 4 Queries, Views, and Rewritings

We consider here (tree node) positions given by Dewey-style strings in  $\mathbb{N}^*$ . If  $t$  is a tree, we denote by  $pos(t)$  the set of all  $t$ 's positions.  $pos(t)$  is prefix-closed and contains  $\epsilon$  as the root position. Given a position  $x \in pos(t)$ , we define  $t_x$  to be the subtree of  $t$  rooted at  $x$ .

**Definition 4.** A pattern is a tree-position pair  $(p, x)$ , where  $x \in pos(p)$ .

Let  $\mathcal{Y}^x = \{(p, x) \mid p \in \mathcal{Y}, x \in pos(p)\}$ .

**Definition 5.** A tree query (TQ)  $Q$  is a subset of  $\mathcal{Y}^x$ .

When a query has only one pattern  $(p, x)$ , we will blur the distinction between  $\{(p, x)\}$  and  $(p, x)$ .

Recall that  $\mathcal{Y}$  is the set of all trees over  $\Sigma$ . We call them *target trees*. Let  $t \in \mathcal{Y}$  be a target tree.

**Definition 6.** The answer to  $Q$  on  $t \in \mathcal{Y}$  is  $ans(Q, t) = \{t_x : (t, x) \in Q\}$ .

For two queries  $Q_1, Q_2$ , we define containment and equivalence as follows.

**Definition 7.**

1.  $Q_1 \sqsubseteq Q_2$  if  $ans(Q_1, t) \subseteq ans(Q_2, t)$  for each  $t \in \mathcal{Y}$ .
2.  $Q_1 \equiv Q_2$  if  $Q_1 \sqsubseteq Q_2$  and  $Q_2 \sqsubseteq Q_1$ .

We show that

**Theorem 2.** If  $Q_1 \subseteq Q_2$  then  $Q_1 \sqsubseteq Q_2$ .

*Proof.* Let  $t' \in ans(Q_1, t)$ . Then there exists  $(t, x) \in Q_1$  such that  $t_x = t'$ . Since  $Q_1 \subseteq Q_2$ , we have  $(t, x) \in Q_2$  and thus  $t_x \in ans(Q_2, t)$ .  $\square$

**Corollary 1.** If  $Q_1 = Q_2$  then  $Q_1 \equiv Q_2$ .

**Regular TQs.**

Let  $\hat{\Sigma} = \{\hat{a} : a \in \Sigma\}$  be an alphabet of marked symbols. Given  $(p, x)$ , consider  $\hat{p}$  on  $\Sigma \cup \hat{\Sigma}$  that is the same as  $p$ , but with the node at position  $x$  being marked by the corresponding symbol in  $\hat{\Sigma}$ . Now a query  $Q$  is a *regular tree query* (RTQ) if set  $\{\hat{p} \mid (p, x) \in Q\}$  is regular, i.e. given by a tree automaton.

*Example 3.* Let us revisit Example 2. Assume that the collection of movie trees is organized into a super-tree  $t$  with a root labeled by  $\mathbf{r}$ , and the movie trees as sub-trees of the root. Now let  $Q$  be a query that asks for all the movie sub-trees having at least one character played by two or more actors. It can be verified that this query can be given by the following automaton.

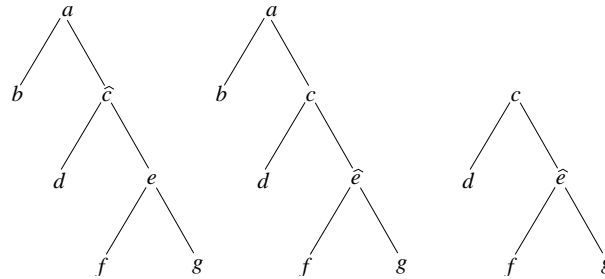
$$\mathcal{A} = (\{s, s_t, s_m, s_c, s_a\}, \Sigma \cup \hat{\Sigma}, \{s_t\}, \Delta)$$

where  $\Delta$  has the following transition rules

$$\begin{aligned}
s^* s_m s^* &\xrightarrow{\mathbf{r}} s_t \\
s^* s_c s^* &\xrightarrow{\hat{\mathbf{m}}} s_m \\
s^* s_a s_a s^* &\xrightarrow{\mathbf{c}} s_c \\
s^* &\xrightarrow{\mathbf{a}} s_a \\
s^* &\xrightarrow{\Sigma} s. \quad \square
\end{aligned}$$

**Views and Rewritings.** First we give a simple example to build up the intuition.

*Example 4.* Consider a view definition containing only the pattern tree given in Figure 2 [left] and a query containing only the pattern tree given in the same figure [middle]. Clearly, the view is useful in answering the query. Intuitively, the pattern tree we need to use to extract the answer to the query using the view is shown in the same figure [right]. This is nothing else but the “rewriting” of the query using the view in this example.



**Fig. 2.** A simple view, a simple query, and the rewriting. Pattern positions are given by  $\hat{\cdot}$ .

We now make precise a rewriting of a query  $Q$  using a view  $V$ .

Recall sets  $\mathcal{T}, \mathcal{T}^x$ . Furthermore, let  $\mathcal{T}^{x:y}$  be the set of all tree-position-position triples  $(p, x, y)$ , such that  $x, y \in \text{pos}(p)$ , and  $x$  is a proper prefix of  $y$ , i.e.  $x$  is a proper ancestor of  $y$ . We restrict  $x$  to be a proper prefix of  $y$  because, as we show in 10.2, the queries and views can be transformed so that their specified positions (markings) never coincide in matching patterns.<sup>3</sup>

Let  $(p, x), (q, y) \in \mathcal{T}^x$ . We define

$$(p, x) \star (q, y) = \begin{cases} (p, xy) & \text{if } p_x = q \text{ and } y \neq \epsilon \\ \text{undefined} & \text{otherwise} \end{cases}$$

<sup>3</sup> 10.2 should be read after being familiarized with the notation introduced in 5.



and

$$(p, x) \star (q, y) = \begin{cases} (p, x, xy) & \text{if } p_x = q \text{ and } y \neq \epsilon \\ \text{undefined} & \text{otherwise.} \end{cases}$$

where  $xy$  is the concatenation of  $x$  and  $y$ .

In simple words, for  $(p, x) \star (q, y)$  and  $(p, x) \star (q, y)$  to be defined, the subtree of  $p$  rooted at  $x$  must be *identical to  $q$*  structure-wise. Also, structure-wise,  $(p, x) \star (q, y)$  and  $(p, x) \star (q, y)$  are the same as  $(p, x)$ .

We note that, when defined,  $(p, x) \star (q, y) \in \mathcal{T}^x$  and  $(p, x) \star (q, y) \in \mathcal{T}^{x,y}$ . Referring to Figure 2 from the previous example, let the view be the tree on the left and the query be the tree in the middle. Then the tree on the right is a rewriting. If we regard the view on the left as  $(p, x)$  and the query in the middle as  $(p, x) \star (q, y)$ , then the rewriting needed to answer the query is  $(q, y)$ , which is indeed the tree on the right. Therefore, finding the rewriting is solving a  $\star$  equation with the rewriting as the unknown.

For two sets  $L_1, L_2 \subseteq \mathcal{T}^x$ , we define  $L_1 \star L_2$  and  $L_1 \star L_2$  in the natural way. Also, we blur the distinction between a set of one element and the element itself. We note that  $L_1 \star L_2 \subseteq \mathcal{T}^x$  and  $L_1 \star L_2 \subseteq \mathcal{T}^{x,y}$ .

We will use  $\star$  to define the rewriting and  $\star$  to reason about some steps in the construction for the rewriting.

In the following we use Greek letters  $v$  and  $\xi$  to denote tree-position pairs.

**Definition 8.** *The maximally contained rewriting (MCR) of  $Q \subseteq \mathcal{T}^x$  using  $V \subseteq \mathcal{T}^x$  is*

$$R = \{\xi \in \mathcal{T}^x : V \star \xi \subseteq Q\}.$$

We also define set  $X$  of “bad” patterns as

$$X = \{\xi \in \mathcal{T}^x : \text{there exists } v \in V \text{ such that } v \star \xi \in Q^c\}$$

where  $Q^c = \mathcal{T}^x \setminus Q$ . Set  $X$  will be crucial in our construction. Observe that  $X$  is not equal to  $\mathcal{T}^x \setminus R$ . The latter also contains patterns that cannot be  $\star$ -ed with  $V$ 's patterns. Now let us define

$$Y = \{\xi \in \mathcal{T}^x : \text{there exists } v \in V \text{ such that } v \star \xi \in Q\}.$$

We have that

**Proposition 2.**  $R = Y \setminus X$ .

*Example 5.* Consider the following queries and views on chain trees, where the positions have been represented as  $\hat{\cdot}$  over the corresponding nodes. The chains should be read backwards to simulate a bottom-up processing.

1. Let  $Q = \{aac\hat{d}e, bbc\hat{d}e, ccc\hat{d}e\}$  and  $V = \{aa\hat{c}de, bb\hat{c}de\}$ . We have  $c\hat{d}e \in R$  because  $V \star c\hat{d}e = \{aac\hat{d}e, bbc\hat{d}e\} \subseteq Q$ . In fact  $c\hat{d}e$  is the only chain in  $R$ .

2. Let  $Q = \{aac\hat{d}e, ccc\hat{d}e\}$  and  $V = \{aa\hat{c}de, bb\hat{c}de\}$ . We have  $c\hat{d}e \notin R$  because  $bbc\hat{d}e \in V \star c\hat{d}e$ , but  $bbc\hat{d}e \notin Q$ . Notice,  $c\hat{d}e$  is in  $X$  and  $Y$ .  $\square$

Let  $T_{V,t}$  be the *materialized answer* (MA) to  $V$  on  $t \in \mathcal{Y}$ , i.e.,  $T_{V,t} = ans(V, t)$ . We answer  $Q$  using  $V$  by computing

$$ans(R, T_{V,t}) = \bigcup_{v \in T_{V,t}} ans(R, v)$$

where  $R$  is the MCR of  $Q$  using  $V$ . We have that

**Theorem 3.**  $ans(R, T_{V,t}) \subseteq ans(Q, t)$  for each  $t \in \mathcal{Y}$ .

*Proof.* From Definition 8 and Theorem 2, we have  $ans(V \star R, t) \subseteq ans(Q, t)$ . Now the claim follows from this equality which we prove next

$$ans(V \star R, t) = ans(R, T_{V,t}).$$

“ $\supseteq$ ”: Let  $z \in ans(R, T_{V,t})$ . Then there exists  $(q, y) \in R$  and  $q \in T_{V,t}$  such that  $q_y = z$ . Since  $q \in T_{V,t}$ , there exists  $(t, x) \in V$  such that  $t_x = q$ . Clearly,  $q_y = z \in ans((t, x) \star (q, y), t)$ . By  $(t, x) \in V$  and  $(q, y) \in R$  we obtain  $z \in ans(V \star R, t)$  as required.

“ $\subseteq$ ”: Let  $z \in ans(V \star R, t)$ . There exists  $(t, xy) = (t, x) \star (q, y) \in V \star R$ , where  $(t, x) \in V$  and  $(q, y) \in R$ , such that  $ans((t, xy), t) = t_{xy} = z$ . Consequently,  $t_x \in T_{V,t}$  and  $z \in ans(R, T_{V,t})$ .  $\square$

## 5 Computing the MCR

The previous section refers to general queries and views. In this section we show that when they are RTQs, we can effectively compute QRV and show that it is an RTQ as well.

First we define the inverse of the  $\star$  operation. Let  $v, v' \in \mathcal{Y}^x$ . We define

$$v \circledast v' = \begin{cases} \xi & \text{if } v' = v \star \xi \\ \text{undefined} & \text{otherwise} \end{cases}$$

This definition is lifted to subsets of  $\mathcal{Y}^x$  in the natural way. Now, it can be verified that

**Proposition 3.**  $X = V \circledast Q^c$  and  $Y = V \circledast Q$ .

Once we obtain  $X$  and  $Y$ , we can obtain  $R$  as  $Y \setminus X$  (see Proposition 2). In the rest of the section, we present an automata-based solution for computing

$$K = J \circledast J'$$

when  $J$  and  $J'$  are RTQs.

Tree-position pair-sets facilitate the definitions of queries and rewritings. However, when we work with automata, we need to talk in terms of languages that the automata recognize. Here, we introduce colors for marking the positions in pattern trees. Colors are similar to Boolean markings of [28] and [23] (or to our previous  $\hat{\phantom{a}}$  marking), however, they make the presentation easier.

There to, we consider the “red” and “blue” alphabets  $\Sigma^{\mathbf{r}} = \{a^{\mathbf{r}} : a \in \Sigma\}$  and  $\Sigma^{\mathbf{b}} = \{a^{\mathbf{b}} : a \in \Sigma\}$  in addition to the alphabet  $\Sigma$ . We refer to the elements of  $\Sigma$  as being “black” and to  $\Sigma$  as the “black” alphabet. We refer to nodes as black, red, or blue nodes if their symbol is black, red, or blue.

We define  $\mathbf{b}$  as a (unary) operator that given  $(p, x) \in \mathcal{T}^{\mathbf{x}}$  returns a tree over  $\Sigma \cup \Sigma^{\mathbf{b}}$  that is isomorphic to  $p$  with the node at position  $x$  colored blue.

Likewise, we define  $\mathbf{r}$  as a (unary) operator that given  $(p, y) \in \mathcal{T}^{\mathbf{y}}$  returns a tree over  $\Sigma \cup \Sigma^{\mathbf{r}}$  that is isomorphic to  $p$  with the node at position  $y$  colored red. [ $\mathcal{T}^{\mathbf{y}}$  is the same as  $\mathcal{T}^{\mathbf{x}}$  and is used for notation parallelism.]

Furthermore, we define  $\mathbf{br}$  as a (unary) operator that given  $(p, x, y) \in \mathcal{T}^{\mathbf{x}, \mathbf{y}}$  returns a tree over  $\Sigma \cup \Sigma^{\mathbf{b}} \cup \Sigma^{\mathbf{r}}$  that is isomorphic to  $p$  with the node at position  $x$  colored blue and the node at the position  $y$  colored red.

Based on these operators, we define

$$\begin{aligned}\mathcal{T}^{\mathbf{b}} &= \{\mathbf{b}(p, x) : (p, x) \in \mathcal{T}^{\mathbf{x}}\} \\ \mathcal{T}^{\mathbf{r}} &= \{\mathbf{r}(p, y) : (p, y) \in \mathcal{T}^{\mathbf{y}}\} \\ \mathcal{T}^{\mathbf{b}, \mathbf{r}} &= \{\mathbf{br}(p, x, y) : (p, x, y) \in \mathcal{T}^{\mathbf{x}, \mathbf{y}}\}.\end{aligned}$$

$\mathcal{T}^{\mathbf{b}}$ ,  $\mathcal{T}^{\mathbf{r}}$ ,  $\mathcal{T}^{\mathbf{b}, \mathbf{r}}$  are *languages of trees* over  $\Sigma \cup \Sigma^{\mathbf{b}}$ ,  $\Sigma \cup \Sigma^{\mathbf{r}}$ ,  $\Sigma \cup \Sigma^{\mathbf{b}} \cup \Sigma^{\mathbf{r}}$ , respectively.

Clearly, there is a one-to-one correspondence between the elements of  $\mathcal{T}^{\mathbf{x}}$ ,  $\mathcal{T}^{\mathbf{y}}$ ,  $\mathcal{T}^{\mathbf{x}, \mathbf{y}}$ , and the elements of  $\mathcal{T}^{\mathbf{b}}$ ,  $\mathcal{T}^{\mathbf{r}}$ ,  $\mathcal{T}^{\mathbf{b}, \mathbf{r}}$ , respectively. Therefore, we will blur the distinction between elements and subsets of  $\mathcal{T}^{\mathbf{x}}$ ,  $\mathcal{T}^{\mathbf{y}}$ ,  $\mathcal{T}^{\mathbf{x}, \mathbf{y}}$ , and elements and subsets of  $\mathcal{T}^{\mathbf{b}}$ ,  $\mathcal{T}^{\mathbf{r}}$ ,  $\mathcal{T}^{\mathbf{b}, \mathbf{r}}$ , respectively.

We will use (sans-serif)  $\mathbf{p}, \mathbf{q}$  to refer to colored patterns of  $\mathcal{T}^{\mathbf{b}}$ ,  $\mathcal{T}^{\mathbf{r}}$ ,  $\mathcal{T}^{\mathbf{b}, \mathbf{r}}$ .

Let  $\mathbf{p} \in \mathcal{T}^{\mathbf{b}}$  and  $\mathbf{q} \in \mathcal{T}^{\mathbf{r}}$  with their corresponding  $(p, x) \in \mathcal{T}^{\mathbf{x}}$  and  $(q, y) \in \mathcal{T}^{\mathbf{y}}$ . We define

$$\begin{aligned}\mathbf{p} \star \mathbf{q} &= \mathbf{r}((p, x) \star (q, y)) \\ \mathbf{p} \star \mathbf{q} &= \mathbf{br}((p, x) \star (q, y)).\end{aligned}$$

We have  $\mathbf{p} \star \mathbf{q} \in \mathcal{T}^{\mathbf{r}}$  and  $\mathbf{p} \star \mathbf{q} \in \mathcal{T}^{\mathbf{b}, \mathbf{r}}$  (when they are defined). We extend  $\star$  and  $\star$  to languages in the natural way.

In particular, from now on, we will consider  $J \subseteq \mathcal{T}^{\mathbf{b}}$ ,  $J' \subseteq \mathcal{T}^{\mathbf{r}}$ ,  $K \subseteq \mathcal{T}^{\mathbf{r}}$ ,  $J \star K \subseteq \mathcal{T}^{\mathbf{r}}$ , and  $J \star K \subseteq \mathcal{T}^{\mathbf{b}, \mathbf{r}}$ .

In order to aid us in the development, we also define

$\Phi^{\mathbf{r}}$  to be the set of all trees having all nodes black, except for the root which is red, and

$\Phi^{\mathbf{b}, \mathbf{r}}$  to be the set of all trees having all nodes black, except for the root which is blue and another node which is red.

It can be verified that the languages we defined so far,  $\mathcal{Y}^r$ ,  $\mathcal{Y}^b$ ,  $\mathcal{Y}^{b,r}$ ,  $\Phi^r$ , and  $\Phi^{b,r}$  are all RTLs. Also observe that from the definition of the  $\clubsuit$  operation,  $K \subseteq \mathcal{Y}^r \setminus \Phi^r$ .

If  $\mathfrak{p}$  is a tree over  $\Sigma \cup \Sigma^r \cup \Sigma^b$ , we denote by  $\mathfrak{p}^{-b}$  the tree over  $\Sigma \cup \Sigma^r$  that is the same as  $\mathfrak{p}$ , but with the blue nodes turned black. For a language  $L$  over  $\Sigma \cup \Sigma^r \cup \Sigma^b$ , we define

$$L^{-b} = \{\mathfrak{p}^{-b} : \mathfrak{p} \in L\}.$$

If  $L$  is an RTL, we can construct an FTA for  $L$  and then an FTA for  $L^{-b}$  by changing all the blue symbols in the transitions of the FTA for  $L$  to black. We define similarly  $\mathfrak{p}^{-r}$  and  $L^{-r}$ .

### 5.1 Auxiliaries

In our construction we will need the following items

$$\begin{aligned} B_L &= \{\mathfrak{p} \in \mathcal{Y}^{b,r} : \mathfrak{p}^{-b} \in L\} \text{ for } L \subseteq \mathcal{Y}^r \\ B'_L &= \{\mathfrak{p} \in \mathcal{Y}^{b,r} : \mathfrak{p}^{-r} \in L\} \text{ for } L \subseteq \mathcal{Y}^b \\ C_L &= \{\mathfrak{p} \in \Phi^{b,r} : \mathfrak{p}^{-b} \in L\} \text{ for } L \subseteq \mathcal{Y}^r \setminus \Phi^r. \end{aligned}$$

These languages are easy to construct when  $L$  is RTL. We show how to do that for  $B_L$ . The other ones are similar. Let  $\mathcal{A} = (S, \Sigma \cup \Sigma^r, F, \Delta)$  be an FTA for  $L$ . We construct FTA  $\mathcal{B} = (S, \Sigma \cup \Sigma^b \cup \Sigma^r, F, \Delta_{\mathcal{B}})$ , where

$$\Delta_{\mathcal{B}} = \Delta \cup \{H \xrightarrow{a^b} s : H \xrightarrow{a} s \text{ in } \Delta \text{ and } a \in \Sigma\}.$$

**Proposition 4.** *Given an FTA for  $L$ , an FTA for  $B_L$  can be constructed in polynomial time. Furthermore, if the FTA for  $L$  is a complete and deterministic FTA (CDFTA), then a CDFTA for  $B_L$  can be constructed in polynomial time as well.*

*Proof.* The statements follow from the observation that  $L(\mathcal{B}) \cap \mathcal{Y}^{b,r} = B_L$ , and the facts that (1)  $\mathcal{B}$  is constructed in polynomial time from  $\mathcal{A}$  and is a CDFTA when  $\mathcal{A}$  is such, (2) an intersection CDFTA is computable in polynomial time when supplied with CDFTAs as input.  $\square$

We also get exactly the same facts for  $B'_L$  and  $C_L$  as those stated for  $B_L$  in the above proposition.

### 5.2 The Algorithm

Here we construct an FTA for  $C_K$  [ $K$ , we are interested in, is  $C_K^{-b}$ , easily computed once we have an automaton for  $C_K$ ].

Consider  $B_J$  and  $B'_{J'}$ . From FTAs for  $B_J$  and  $B'_{J'}$ , we construct a FTA for  $B_J \cap B'_{J'}$ . Let  $\mathcal{D} = (S, \Sigma \cup \Sigma^b \cup \Sigma^r, F, \Delta)$  be this FTA. We also transform it to be targeted. Observe that  $L(\mathcal{D}) \subseteq \mathcal{Y}^{b,r}$ .

Now, we construct FTA  $\mathcal{E} = (S, \Sigma \cup \Sigma^b \cup \Sigma^r, F_{\mathcal{E}}, \Delta)$ , where

$$F_{\mathcal{E}} = \{s \in S : \text{there exists } H \xrightarrow{a^b} s \text{ in } \Delta\}.$$

Clearly,  $L(\mathcal{E}) \subseteq \Phi^{b,r}$ . This is true because  $\mathcal{D}$  is targeted. We have

**Theorem 4.**  $L(\mathcal{E}) = C_K$ .

*Proof.* “ $\subseteq$ ”: Let  $q \in L(\mathcal{E}) \subseteq \Phi^{\mathbf{b},\mathbf{r}}$ . By the construction of  $\mathcal{E}$ , there exists  $p \in B_J \cap B'_{J'} \subseteq \Upsilon^{\mathbf{b},\mathbf{r}}$ , with the blue node at a position  $x$ , such that  $p_x = q$ .

Since  $p \in B_J \cap B'_{J'}$ ,  $p^{-\mathbf{r}} \in J$  and  $p^{-\mathbf{b}} \in J'$ . Also, since  $p \in \Upsilon^{\mathbf{b},\mathbf{r}}$ ,  $(p^{-\mathbf{b}})_x \in \Upsilon^{\mathbf{r}} \setminus \Phi^{\mathbf{r}}$ . We have  $p^{-\mathbf{r}} \star (p^{-\mathbf{b}})_x = p^{-\mathbf{b}}$ , i.e.  $(p^{-\mathbf{b}})_x \in J \star J' = K$ , which means  $p_x = q \in C_K$ .

“ $\supseteq$ ”: Let  $q \in C_K \subseteq \Phi^{\mathbf{b},\mathbf{r}}$ . We have  $q^{-\mathbf{b}} \in K \subseteq \Upsilon^{\mathbf{r}} \setminus \Phi^{\mathbf{r}}$ . By the definition of  $K$ , there exists  $p \in J$  such that  $p \star q^{-\mathbf{b}}$  is defined and  $p \star q^{-\mathbf{b}} \in J'$ . This, coupled with the fact that  $q^{-\mathbf{b}} \in \Upsilon^{\mathbf{r}} \setminus \Phi^{\mathbf{r}}$ , implies  $p \star q^{-\mathbf{b}} \in B_J$ . By the definition of “ $\star$ ” and  $B'_{J'}$ ,  $p \star q^{-\mathbf{b}} \in B'_{J'}$ , too, therefore,  $p \star q^{-\mathbf{b}} \in B_J \cap B'_{J'}$ , i.e.,  $p \star q^{-\mathbf{b}}$  is accepted by  $\mathcal{D}$ . By the construction of  $\mathcal{E}$ , we have that  $q$  is accepted by  $\mathcal{E}$ .  $\square$

### 5.3 Complexity

Let  $J$  and  $J'$  be given by non-deterministic FTAs. It can be verified that

**Proposition 5.**  $K = J \star J'$  can be computed in polynomial time.

Therefore,  $Y = V \star Q$  can be computed in polynomial time. However, we also need to compute  $X = V \star Q^c$ , and then  $X^c$ , in order to compute  $R$  as  $Y \setminus X$ . If we are not careful, we can incur a double-exponential penalty. Here we show that it can be done instead in single-exponential time.

For this let us refer to the steps of computing  $C_K$ . Suppose  $J$  and  $J'$  are given by CDFTAs. By Proposition 4 we can obtain CDFTAs for  $B_J$  and  $B'_{J'}$  in polynomial time. From these CDFTAs we obtain a CDFTA  $\mathcal{D}$  for  $B_J \cap B'_{J'}$  in polynomial time as well. Automaton  $\mathcal{E}$  (representing  $C_K$ ) is the same as  $\mathcal{D}$ , but with a different set of final states, so  $\mathcal{E}$  is a CDFTA as well.

Since  $\mathcal{E}$  is a CDFTA we can compute  $(C_K)^c$  from it in polynomial time.

Now, we have that for  $X$  and  $Y$

1.  $C_X, C_Y \subseteq \Phi^{\mathbf{b},\mathbf{r}}$
2.  $Y \setminus X = (C_Y \setminus C_X)^{-\mathbf{b}} = (C_Y \cap (C_X)^c)^{-\mathbf{b}}$ .

Clearly,  $\cap$  and  $-\mathbf{b}$  are polynomial.

Let a query  $Q$  and a view  $V$  be specified as non-deterministic FTAs. We have the following theorems.

**Theorem 5.** The MCR  $R$  of  $Q$  using  $V$  can be computed in exponential time.

*Proof.* The proof follows from the above reasoning and the fact that to construct CDFTAs for  $Q$  (and then  $Q^c$ ), and  $V$  takes exponential time.  $\square$

**Theorem 6.** Computing the MCR of  $Q$  using  $V$  is EXPTIME-hard.

*Proof.* We present a reduction from the problem of inclusion for RTLs (represented by non-deterministic FTAs with horizontal languages represented by NFAs) which is EXPTIME-complete. Let  $L_1$  and  $L_2$  be such RTLs over some alphabet  $\Sigma$ . Let  $\# \notin \Sigma$ . We construct  $V$  from  $L_1$  by adding to the rightmost

leaf of trees in  $L_1$  a subtree  $\#^b(\#)$ . Similarly, we construct  $Q$  from  $L_2$  by adding to the rightmost leaf of trees in  $L_2$  a subtree  $\#(\#^r)$ . Now, we have that  $L_1 \subseteq L_2$  if and only if the MCR of  $Q$  using  $V$  contains the single node pattern  $\#^r$ . The latter can be checked in time polynomial in the size of an FTA for the rewriting.

In conclusion, if computing the MCR of  $Q$  using  $V$  were not EXPTIME-hard, then we would be able to decide the language inclusion for RTLs in less than EXPTIME, which is not possible in the worst case.  $\square$

## 6 $k$ -ary Queries

So far, our attention has been focused on unary queries. In this section, we consider  $k$ -ary queries, for  $k \geq 1$ . Miklau and Suciu [18] motivate the study of containment for  $k$ -ary tree pattern queries by observing that the results can be applied in the context of optimizing the FOR clause of XQuery expressions, where multiple variables can be defined using XPath variables. The idea is that they can be captured in a single tree pattern<sup>4</sup> whose distinguished variables correspond to the variables bound in the FOR clause. Inspired by this, we consider tree-position vector queries whose vector contains  $k$  positions, for some arbitrary, but *fixed*  $k \geq 1$ . Such queries return a tuple (forest) of sub-trees from a target tree.

On the theoretical side, it has been shown that  $k$ -ary queries can be encoded by unary queries (see [25] and [3]). However, this is done by going through MSO formulas. Therefore, going from a  $k$ -ary query represented by an automaton to an encoding by unary queries and then back to (combination of) automata would incur non elementary complexity. Even if the complexity were not so, we still would be left with a combination of unary queries for which we needed to compute QRV in terms of another combination of unary queries which corresponds to the view. This is something that does not follow from the rewriting mechanism we developed in the previous part of the paper. The latter comment also applies to other  $k$ -ary query formalisms, defined in terms of compositions of unary queries such as the formalism introduced in [13].

**Definition 9.** A  $k$ -pattern is a tree-position vector pair  $(p, [x_1, \dots, x_k])$ , where

1.  $x_1, \dots, x_k \in \text{pos}(p)$
2.  $x_1 < \dots < x_k$
3.  $\nexists i, j \in [1, k]$ , such that  $x_i$  is a prefix of  $x_j$ .

We also write a pattern as  $(p, \mathbf{x})$ , where  $\mathbf{x} = [x_1, \dots, x_k]$ . A position vector possessing the three properties in above definition is called *proper*. Whenever we refer to a position vector, we assume it to be proper. We let  $\mathcal{T}^{k \times \mathbf{x}}$  denote the set of all  $k$ -patterns.

Given a tree  $t$  and  $\mathbf{x} = [x_1, \dots, x_k]$ , we define  $t_{\mathbf{x}} = (t_{x_1}, \dots, t_{x_k})$  as the tuple (forest) of sub-trees of  $t$  rooted at the positions in  $\mathbf{x}$ . We call a forest of  $k$  trees a  $k$ -forest, and denote by  $\Xi^k$  the set of all  $k$ -forests.

<sup>4</sup> Tree patterns considered in [18] have child and descendant edges, wildcards, and no order.

**Definition 10.** A  $k$ -tree query ( $k$ -TQ) is a subset of  $\mathcal{Y}^{k\mathbf{x}}$ . If  $\{\hat{p} : (p, \mathbf{x}) \in Q\}$ <sup>5</sup> is in addition regular, we say that  $Q$  is a regular tree query (RTQ).

When a query has only one pattern  $(p, \mathbf{x})$ , we blur the distinction between  $\{(p, \mathbf{x})\}$  and  $(p, \mathbf{x})$ .

**Definition 11.** The answer to a  $k$ -query  $Q$  on  $t$  is  $\text{ans}(Q, t) = \{t_{\mathbf{x}} : (t, \mathbf{x}) \in Q\}$ .

As for unary queries, it can be shown that for  $k$ -tree queries, set containment implies query containment.

*Example 6.* Let us recall the movie example in Section 4. Now let  $Q$  be a query given by  $\mathcal{A} = (\{s, s_t, s_m, s_c, s_a\}, \Sigma \cup \hat{\Sigma}, \{s_t\}, \Delta)$ , where  $\Delta$  has the following transition rules

$$\begin{aligned} s^* s_m s^* &\xrightarrow{\mathbf{r}} s_t \\ s^* s_c s^* &\xrightarrow{\mathbf{m}} s_m \\ s^* s_a s^* s_a s^* &\xrightarrow{\mathbf{c}} s_c \\ s^* &\xrightarrow{\hat{\mathbf{a}}} s_a \\ s^* &\xrightarrow{\Sigma} s. \end{aligned}$$

Clearly,  $Q \subseteq \mathcal{Y}^{2\mathbf{x}}$ , and  $\text{ans}(Q, t)$  consists of all 2-forests of (sub)tree pairs for actors who have played the same character together in some movie.  $\square$

## 7 $k$ -ary Views and Rewritings

We define a  $k, m\mathbf{y}$ -forest to be a tuple  $((q_1, \mathbf{y}_1), \dots, (q_k, \mathbf{y}_k))$  of patterns, such that  $|\mathbf{y}_1| + \dots + |\mathbf{y}_k| = m$ , where  $|\mathbf{y}_i|$ , for  $i \in [1, k]$ , is the dimension of vector  $\mathbf{y}_i$ . We let  $\Xi^{k, m\mathbf{y}}$  be the set of all  $k, m\mathbf{y}$ -forests.

Furthermore, we define a  $k\mathbf{x}, m\mathbf{y}$ -pattern to be a tuple of the form

$$(p, [(x_1, \mathbf{y}_1), \dots, (x_k, \mathbf{y}_k)])$$

where  $|\mathbf{y}_1| + \dots + |\mathbf{y}_k| = m$ , and  $x_i$  is a proper<sup>6</sup> prefix of positions in  $\mathbf{y}_i$ , for  $i \in [1, k]$ . We let  $\mathcal{Y}^{k\mathbf{x}, m\mathbf{y}}$  be the set of all  $k\mathbf{x}, m\mathbf{y}$ -patterns.

Now, let  $(p, [x_1, \dots, x_k]) \in \mathcal{Y}^{k\mathbf{x}}$  and  $((q_1, \mathbf{y}_1), \dots, (q_k, \mathbf{y}_k)) \in \Xi^{k, m\mathbf{y}}$ . We define

$$\begin{aligned} (p, [x_1, \dots, x_k]) \star ((q_1, \mathbf{y}_1), \dots, (q_k, \mathbf{y}_k)) = \\ \begin{cases} (p, [x_1\mathbf{y}_1, \dots, x_k\mathbf{y}_k]) & \text{if } p_{x_i} = q_i \text{ and } y_i \neq \epsilon, \forall i \in [1, k] \\ \text{undefined} & \text{otherwise} \end{cases} \end{aligned}$$

<sup>5</sup> Similarly as for unary queries,  $\hat{p}$  denotes the pattern on  $\Sigma \cup \hat{\Sigma}$  that is the same as  $p$ , but with the nodes at positions of  $\mathbf{x}$  being marked by the corresponding symbols in  $\hat{\Sigma}$ .

<sup>6</sup> By assuming properly transformed query and view (as in 10.2), we do not need to consider the case when some  $x_i$  coincides with some position in  $\mathbf{y}_i$ .

and

$$(p, [x_1, \dots, x_k]) \star ((q_1, \mathbf{y}_1), \dots, (q_k, \mathbf{y}_k)) = \begin{cases} (p, [(x_1, x_1 \mathbf{y}_1), \dots, (x_k, x_k \mathbf{y}_k)]) & \text{if } p_{x_i} = q_i \text{ and } y_i \neq \epsilon, \forall i \in [1, k] \\ \text{undefined} & \text{otherwise} \end{cases}$$

where  $x_i \mathbf{y}_i$ , for  $i \in [1, k]$ , is the tuple obtained from  $\mathbf{y}_i$  by prepending  $x_i$  to each of its positions.

We note that, when defined

$$\begin{aligned} (p, [x_1, \dots, x_k]) \star ((q_1, \mathbf{y}_1), \dots, (q_k, \mathbf{y}_k)) &\in \mathcal{Y}^{m\mathbf{y}} \\ (p, [x_1, \dots, x_k]) \star ((q_1, \mathbf{y}_1), \dots, (q_k, \mathbf{y}_k)) &\in \mathcal{Y}^{k\mathbf{x}, m\mathbf{y}}. \end{aligned}$$

[For ease of exposition we are using  $\mathcal{Y}^{m\mathbf{y}}$  instead of  $\mathcal{Y}^{m\mathbf{x}}$ .]

For  $L_1 \subseteq \mathcal{Y}^{k\mathbf{x}}$ ,  $L_2 \subseteq \Xi^{k, m\mathbf{y}}$  we define  $L_1 \star L_2$  and  $L_1 \star L_2$  in the natural way. Again we blur the distinction between a set of one element and the element itself.

**Definition 12.** *The maximally contained rewriting (MCR) of  $Q \subseteq \mathcal{Y}^{m\mathbf{y}}$  using  $V \subseteq \mathcal{Y}^{k\mathbf{x}}$  is*

$$R = \{\xi \in \Xi^{k, m\mathbf{y}} : V \star \xi \subseteq Q\}.$$

( $\xi$  being in  $\Xi^{k, m\mathbf{y}}$  is of the form  $((q_1, \mathbf{y}_1), \dots, (q_k, \mathbf{y}_k))$ .)

We also define set  $X$  of “bad” patterns as

$$X = \{\xi \in \Xi^{k, m\mathbf{y}} : \text{there exists } v \in V \text{ s.t. } v \star \xi \in Q^c\}$$

where  $Q^c = \mathcal{Y}^{m\mathbf{y}} \setminus Q$ . ( $v$  being in  $\mathcal{Y}^{k\mathbf{x}}$  is of the form  $(p, [x_1, \dots, x_k])$ .) Now we define

$$Y = \{\xi \in \Xi^{k, m\mathbf{y}} : \text{there exists } v \in V \text{ s.t. } v \star \xi \in Q\}.$$

Notice that  $R$ ,  $X$ , and  $Y$  are subsets of  $\Xi^{k, m\mathbf{y}}$ . We have that

**Proposition 6.**  $R = Y \setminus X$ .

**Dummy roots.** We turn forests into trees by introducing a “dummy” root labeled by a special symbol  $\diamond \notin \Sigma$ .

Given  $(t_1, \dots, t_k) \in \Xi^k$ , we define

$$\diamond(t_1, \dots, t_k)$$

to be the tree obtained by making the roots of  $t_1, \dots, t_k$  to be children of a dummy new root labeled by  $\diamond$ .

Likewise, given  $((q_1, \mathbf{y}_1), \dots, (q_k, \mathbf{y}_k)) \in \Xi^{k, m\mathbf{y}}$ , we define

$$\diamond((q_1, \mathbf{y}_1), \dots, (q_k, \mathbf{y}_k))$$

to be the tree obtained by making the roots of  $q_1, \dots, q_k$  to be children of a dummy new root labeled by  $\diamond$ .



If  $L$  is a subset of  $\Xi^k$  or  $\Xi^{k,m\mathbf{y}}$ , we define  $\diamond L$  in the natural way by turning each forest in  $L$  into a tree as described.

Let  $F \subseteq \Xi^k$  and  $L \subseteq \Xi^{k,m\mathbf{y}}$ . We define

$$ans(L, F) = \bigcup_{v \in F} ans(\diamond L, \diamond v).$$

Let  $F_{V,t} = ans(V, t) \subseteq \Xi^k$  be the materialization of  $V$  on  $t \in \mathcal{T}$ . We answer  $Q$  using  $V$  by  $ans(R, F_{V,t})$ . Now, by generalizing the reasoning in the proof of Theorem 3 it can be verified that

**Theorem 7.**  $ans(R, F_{V,t}) \subseteq ans(Q, t)$  for each  $t \in \mathcal{T}$ .

## 8 Computing the MCR of an $m$ -ary Query using a $k$ -ary View

Analogous to the case of unary queries, we first define the inverse of the  $\star$  operation.

Let  $v \in \mathcal{Y}^{k\mathbf{x}}$  and  $v' \in \mathcal{Y}^{m\mathbf{x}}$ . We define

$$v \circledast v' = \begin{cases} \xi & \text{if } v' = v \star \xi \\ \text{undefined} & \text{otherwise} \end{cases}$$

This definition is lifted to subsets of  $\mathcal{Y}^{\mathbf{x}}$  in the natural way. Now, it can be verified that as in the case of unary queries

**Proposition 7.**  $X = V \circledast Q^c$  and  $Y = V \circledast Q$ .

Once we obtain  $X$  and  $Y$ , we can obtain  $R$  as  $Y \setminus X$ . In the rest of the section, we present an automata-based solution for computing

$$K = J \circledast J'$$

when  $J$  and  $J'$  are RTQs.

Here we again use colors in order to work with languages and automata. We define  $\mathbf{b}$  as a (unary) operator that given  $(p, \mathbf{x}) \in \mathcal{Y}^{k\mathbf{x}}$  returns a tree over  $\Sigma \cup \Sigma^{\mathbf{b}}$  that is isomorphic to  $p$  with the nodes at the positions of  $\mathbf{x}$  colored blue.

Likewise, we define  $\mathbf{r}$  as a (unary) operator that given  $(p, \mathbf{y}) \in \mathcal{Y}^{m\mathbf{y}}$  returns a tree over  $\Sigma \cup \Sigma^{\mathbf{r}}$  that is isomorphic to  $p$  with the nodes at the positions of  $\mathbf{y}$  colored red.

Furthermore, we define  $\mathbf{br}$  as a (unary) operator that given  $(p, (x_1, \mathbf{y}_1), \dots, (x_k, \mathbf{y}_k)) \in \mathcal{Y}^{k\mathbf{x}, m\mathbf{y}}$  returns a tree over  $\Sigma \cup \Sigma^{\mathbf{b}} \cup \Sigma^{\mathbf{r}}$  that is isomorphic to  $p$  with the nodes at positions  $x_1, \dots, x_k$  colored blue and the nodes at the positions of  $\mathbf{y}_1, \dots, \mathbf{y}_k$  colored red.

Based on these operators, we define

$$\begin{aligned}
\mathcal{Y}^{k\mathbf{b}} &= \{\mathbf{b}(p, \mathbf{x}) : (p, \mathbf{x}) \in \mathcal{Y}^{k\mathbf{x}}\} \\
\mathcal{Y}^{m\mathbf{r}} &= \{\mathbf{r}(p, \mathbf{y}) : (p, \mathbf{y}) \in \mathcal{Y}^{m\mathbf{y}}\} \\
\mathcal{Y}^{k\mathbf{b}, m\mathbf{r}} &= \{\mathbf{br}(p, (x_1, \mathbf{y}_1), \dots, (x_k, \mathbf{y}_k)) : \\
&\quad (p, (x_1, \mathbf{y}_1), \dots, (x_k, \mathbf{y}_k)) \in \mathcal{Y}^{k\mathbf{x}, m\mathbf{y}}\} \\
\Xi^{k, m\mathbf{r}} &= \{\mathbf{r}(q_1, \mathbf{y}_1), \dots, \mathbf{r}(q_k, \mathbf{y}_k) : \\
&\quad ((q_1, \mathbf{y}_1), \dots, (q_k, \mathbf{y}_k)) \in \Xi^{k, m\mathbf{y}}\}.
\end{aligned}$$

Clearly, there is a one-to-one correspondence between the elements of  $\mathcal{Y}^{k\mathbf{x}}$ ,  $\mathcal{Y}^{m\mathbf{y}}$ ,  $\mathcal{Y}^{k\mathbf{x}, m\mathbf{y}}$ ,  $\Xi^{k, m\mathbf{y}}$  and the elements of  $\mathcal{Y}^{k\mathbf{b}}$ ,  $\mathcal{Y}^{m\mathbf{r}}$ ,  $\mathcal{Y}^{k\mathbf{b}, m\mathbf{r}}$ ,  $\Xi^{k, m\mathbf{r}}$ , respectively. Therefore, we will blur the distinction between elements and subsets of  $\mathcal{Y}^{k\mathbf{x}}$ ,  $\mathcal{Y}^{m\mathbf{y}}$ ,  $\mathcal{Y}^{k\mathbf{x}, m\mathbf{y}}$ ,  $\Xi^{k, m\mathbf{y}}$  and elements and subsets of  $\mathcal{Y}^{k\mathbf{b}}$ ,  $\mathcal{Y}^{m\mathbf{r}}$ ,  $\mathcal{Y}^{k\mathbf{b}, m\mathbf{r}}$ ,  $\Xi^{k, m\mathbf{r}}$ , respectively.

We will use (sans-serif)  $\mathbf{p}, \mathbf{q}$  to refer to colored patterns of  $\mathcal{Y}^{k\mathbf{b}}$ ,  $\mathcal{Y}^{m\mathbf{r}}$ ,  $\mathcal{Y}^{k\mathbf{b}, m\mathbf{r}}$ , and  $\mathbf{q} = (\mathbf{q}_1, \dots, \mathbf{q}_k)$  to refer to colored forests of  $\Xi^{k, m\mathbf{r}}$ .

Let  $\mathbf{p} \in \mathcal{Y}^{k\mathbf{b}}$  and  $\mathbf{q} = (\mathbf{q}_1, \dots, \mathbf{q}_k) \in \Xi^{k, m\mathbf{r}}$  with their corresponding  $(p, \mathbf{x}) \in \mathcal{Y}^{k\mathbf{x}}$  and  $((q_1, \mathbf{y}_1), \dots, (q_k, \mathbf{y}_k)) \in \Xi^{k, m\mathbf{y}}$ . We define

$$\begin{aligned}
\mathbf{p} \star \mathbf{q} &= \mathbf{r}((p, \mathbf{x}) \star ((q_1, \mathbf{y}_1), \dots, (q_k, \mathbf{y}_k))) \\
\mathbf{p} \star \mathbf{q} &= \mathbf{br}((p, \mathbf{x}) \star ((q_1, \mathbf{y}_1), \dots, (q_k, \mathbf{y}_k))).
\end{aligned}$$

We extend  $\star$  and  $\star$  to languages and forests in the natural way.

In particular, from now on, we will consider  $J \subseteq \mathcal{Y}^{k\mathbf{b}}$ ,  $J' \subseteq \mathcal{Y}^{m\mathbf{r}}$ ,  $K \subseteq \Xi^{k, m\mathbf{r}}$ ,  $J \star K \subseteq \mathcal{Y}^{m\mathbf{r}}$ , and  $J \star K \subseteq \mathcal{Y}^{k\mathbf{b}, m\mathbf{r}}$ .

We further define a few more notions. Let

$\Phi^{k, m\mathbf{r}}$  be the subset of  $\Xi^{k, m\mathbf{r}}$  forests having *at least one* tree that has its root red.

$\Phi^{k\mathbf{b}, m\mathbf{r}}$  be the set of all  $k$ -forests over  $\Sigma \cup \Sigma^{\mathbf{b}} \cup \Sigma^{\mathbf{r}}$ , such that: (1) the roots of all patterns in the forest are blue, (2) *exactly*  $m$  nodes in total are red, and (3) all the other nodes are black.

Observe that from the definition of the  $\star$  operation,  $K \subseteq \Xi^{k, m\mathbf{r}} \setminus \Phi^{k, m\mathbf{r}}$ .

We define the operators  $(\cdot)^{-\mathbf{b}}$  and  $(\cdot)^{-\mathbf{r}}$ , for patterns, sets of patterns, forests, and sets of forests, over  $\Sigma \cup \Sigma^{\mathbf{b}} \cup \Sigma^{\mathbf{r}}$ , in a manner similar to Section 5.

## 8.1 Auxiliaries

In our construction we will need the following items

$$\begin{aligned}
B_L &= \{\mathbf{p} \in \mathcal{Y}^{k\mathbf{b}, m\mathbf{r}} : \mathbf{p}^{-\mathbf{b}} \in L\} \text{ for } L \subseteq \mathcal{Y}^{m\mathbf{r}} \\
B'_L &= \{\mathbf{p} \in \mathcal{Y}^{k\mathbf{b}, m\mathbf{r}} : \mathbf{p}^{-\mathbf{r}} \in L\} \text{ for } L \subseteq \mathcal{Y}^{k\mathbf{b}} \\
C_L &= \{\mathbf{q} \in \Phi^{k\mathbf{b}, m\mathbf{r}} : \mathbf{q}^{-\mathbf{b}} \in L\} \text{ for } L \subseteq \Xi^{k, m\mathbf{r}} \setminus \Phi^{k, m\mathbf{r}}.
\end{aligned}$$

In plain language  $C_L$  is the set of all forests in  $\Phi^{kb,mr}$  obtained from the forests of  $L$  after making the root of their trees blue.

The above languages are easy to construct when  $L$  ( $\diamond L$  in the third definition) is RTL. For instance, for  $B_L$ , we first construct an automaton  $\mathcal{B}$  exactly as in Subsection 5.1, and then compute  $B_L$  as  $L(\mathcal{B}) \cap \Upsilon^{kb,mr}$ . Similar statements as those made in Proposition 4 can be made here as well.

## 8.2 The Algorithm

In order to be compatible with the framework of FTAs (which work on trees, not forests), here we construct an CDFTA for  $\diamond C_K$ . Set  $\diamond K$  (we are interested in) is easily computed once we have an automaton for  $\diamond C_K$ .

Consider  $B_J$  and  $B'_{J'}$ . From CDFTAs for  $B_J$  and  $B'_{J'}$ , we construct a CDFTA for  $B_J \cap B'_{J'}$ . Let  $\mathcal{D} = (S_{\mathcal{D}}, \Sigma \cup \Sigma^b \cup \Sigma^r, F_{\mathcal{D}}, \Delta_{\mathcal{D}})$  be this CDFTA. We also transform it to be targeted. Observe that  $L(\mathcal{D}) \subseteq \Upsilon^{kb,mr}$ . Let

$$S^b = \{s \in S_{\mathcal{D}} : \text{there exists } H \xrightarrow{a^b} s \text{ in } \Delta_{\mathcal{D}}\}.$$

We denote  $(S^b)^k$  by  $S^{kb}$ . Next, we construct FTA  $\mathcal{E} = (S_{\mathcal{E}}, \Sigma \cup \Sigma^b \cup \Sigma^r \cup \{\diamond\}, F_{\mathcal{E}}, \Delta_{\mathcal{E}})$  where

$$\begin{aligned} S_{\mathcal{E}} &= S_{\mathcal{D}} \cup \{s_{\text{final}}, s_{\text{garbage}}\} \\ F_{\mathcal{E}} &= \{s_{\text{final}}\} \\ \Delta_{\mathcal{E}} &= \Delta_{\mathcal{D}} \cup \{S^{kb} \xrightarrow{\diamond} s_{\text{final}}\} \cup \{S_{\mathcal{E}}^* \setminus S^{kb} \xrightarrow{\diamond} s_{\text{garbage}}\} \cup \\ &\quad \{(s_{\text{garbage}})^* \xrightarrow{a} s_{\text{garbage}} : a \in \Sigma \cup \Sigma^b \cup \Sigma^r\}. \end{aligned}$$

Since  $k$  is fixed, an automaton for  $S_{\mathcal{E}}^* \setminus S^{kb}$  can be constructed in polynomial time in the size of  $S^{kb}$  and  $S_{\mathcal{E}}$ . We can show that the FTA  $\mathcal{E}$  constructed above is a CDFTA.

Clearly,  $L(\mathcal{E}) \subseteq \diamond \Phi^{kr,mr}$ . This is true because  $\mathcal{D}$  is targeted. Now, by generalizing the reasoning in the proof of Theorem 4, we can show that

**Theorem 8.**  $L(\mathcal{E}) = \diamond C_K$ . Furthermore,  $\mathcal{E}$  is complete and deterministic.

*Proof.* “ $\subseteq$ .” Let  $\diamond \mathbf{q} \in L(\mathcal{E})$ , where  $\mathbf{q} \in \Phi^{kb,mr}$ . By the construction of  $\mathcal{E}$ , there exists  $\mathbf{p} \in B_J \cap B'_{J'} \subseteq \Upsilon^{kb,mr}$ , with blue node position-tuple  $\mathbf{x}$ , such that  $\mathbf{p}_{\mathbf{x}} = \mathbf{q}$ .

Since  $\mathbf{p} \in B_J \cap B'_{J'}$ ,  $\mathbf{p}^{-r} \in J$  and  $\mathbf{p}^{-b} \in J'$ . Also, observe that  $\mathbf{p}^{-r} \star (\mathbf{p}^{-b})_{\mathbf{x}} = \mathbf{p}^{-b}$ . These imply that  $(\mathbf{p}^{-b})_{\mathbf{x}} \in J \star J' = K$ , which means  $\mathbf{p}_{\mathbf{x}} = \mathbf{q} \in C_K$ , i.e.  $\diamond \mathbf{q} \in \diamond C_K$ .

“ $\supseteq$ .” Let  $\diamond \mathbf{q} \in \diamond C_K \subseteq \diamond \Phi^{kb,mr}$ , where  $\mathbf{q} \in C_K \subseteq \Phi^{kb,mr}$ , i.e.  $\mathbf{q}^{-b} \in K$ . By the definition of  $K$ , there exists  $\mathbf{p} \in J$  such that  $\mathbf{p} \star \mathbf{q}^{-b}$  is defined and in  $J'$ . This implies  $\mathbf{p} \star \mathbf{q}^{-b} \in B_J$ . By the definition of “ $\star$ ” and  $B'_{J'}$ ,  $\mathbf{p} \star \mathbf{q}^{-b} \in B'_{J'}$ , too, therefore,  $\mathbf{p} \star \mathbf{q}^{-b} \in B_J \cap B'_{J'}$ , i.e.  $\mathbf{p} \star \mathbf{q}^{-b}$  is accepted by  $\mathcal{D}$ . Now, by the construction of  $\mathcal{E}$ , we have  $\diamond \mathbf{q}$  is accepted by  $\mathcal{E}$ .  $\square$

**Complexity.** It can be verified that, similarly as in the case of unary queries (Subsection 5.3), computing  $R$  can be done in singly-exponential time when starting with non-deterministic FTAs for  $Q$  and  $V$ . By Theorem 6, which applies to the case of  $k = m = 1$ , it follows that our constructions are optimal.

## 9 Conclusions

In this paper, we studied the problem of rewriting queries using views for XML data, focusing on tree-selecting queries equal in power to MSO. Colored tree languages and automata provided us a critical ground for implementing rewriting queries using views. We developed a singly-exponential algorithm for deriving the maximally contained rewriting of a query using a view, and then extended it to an algorithm for the general case where the query can be an  $m$ -ary query and the view can be a  $k$ -ary view. The latter class of queries is useful in optimizing the FOR clause of XQuery expressions. We showed the problem is EXPTIME-hard, thus showing our algorithms are optimal. Allowing our query formalism to also involve selection, join and restructuring is an avenue worth exploring. Extending the reasoning about QRV for this wider class is an interesting open problem.

## References

1. L. Afanasiev, P. Blackburn, I. Dimitriou, B. Gaiffe, E. Goris, M. Marx, and M. de Rijke. PDL for ordered trees. *Journal of Applied Non-Classical Logics*, 15(2):115–135, 2005.
2. F. N. Afrati, R. Chirkova, M. Gergatsoulis, B. Kimelfeld, V. Pavlaki, and Y. Sagiv. On rewriting XPath queries using views. In *EDBT*, 2009.
3. M. Arenas, P. Barceló, and L. Libkin. Combining temporal logics for querying XML documents. In *ICDT*, 2007.
4. A. Arion, V. Benzaken, I. Manolescu, and Y. Papakonstantinou. Structured materialized views for XML queries. In *VLDB*, 2007.
5. A. Balmin, F. Özcan, K. S. Beyer, R. Cochrane, and H. Pirahesh. A framework for using materialized XPath views in XML query processing. In *VLDB*, 2004.
6. D. Calvanese, G. D. Giacomo, M. Lenzerini, and M. Y. Vardi. An automata-theoretic approach to regular XPath. In *DBPL*, 2009.
7. D. Calvanese, G. D. Giacomo, M. Lenzerini, and M. Y. Vardi. Node selection query languages for trees. In *AAAI*, 2010.
8. B. Cautis, A. Deutsch, and N. Onose. XPath rewriting using multiple views: Achieving completeness and efficiency. In *WebDB*, 2008.
9. P. Claves, D. Jansen, S. J. Holtrup, M. Mohr, A. Reis, M. Schatz, and I. Thesing. LETHAL: Library for working with finite tree and hedge automata. Available on: <http://lethal.sf.net>, 2009.
10. H. Comon, M. Dauchet, R. Gilleron, C. Löding, F. Jacquemard, D. Lugiez, S. Tison, and M. Tommasi. Tree automata techniques and applications, 2007.
11. W. Fan, F. Geerts, X. Jia, and A. Kementsietsidis. Rewriting regular XPath queries on XML views. In *ICDE*, 2007.

12. E. Filiot. Ranked and unranked tree automata libraries (grappa). Available on: <http://www.grappa.univ-lille3.fr/~filiot/tata>.
13. E. Filiot, J. Niehren, J.-M. Talbot, and S. Tison. Composing monadic queries in trees. In *PLAN-X*, 2006.
14. M. Frick, M. Grohe, and C. Koch. Query evaluation on compressed trees (extended abstract). In *LICS*, 2003.
15. J. E. Hopcroft and J. D. Ullman. *Introduction to Automata Theory, Languages and Computation*. Addison-Wesley, 1979.
16. L. V. S. Lakshmanan, H. Wang, and Z. J. Zhao. Answering tree pattern queries using views. In *VLDB*, 2006.
17. L. Libkin and C. Sirangelo. Reasoning about XML with temporal logics and automata. *J. Applied Logic*, 8(2):210–232, 2010.
18. G. Miklau and D. Suciu. Containment and equivalence for a fragment of XPath. *JACM*, 51(1):2–45, 2004.
19. F. Neven. *Design and Analysis of Query Languages for Structured Documents—A Formal and Logical Approach*. PhD thesis. Limburgs Universitair Centrum, 1999.
20. F. Neven and T. Schwentick. Expressive and efficient pattern languages for tree-structured data. In *PODS*, 2000.
21. F. Neven and T. Schwentick. Query automata over finite trees. *TCS*, 275(1-2):633–674, 2002.
22. F. Neven and T. Schwentick. On the complexity of XPath containment in the presence of disjunction, DTDs, and variables. *Logical Methods in Computer Science*, 2(3), 2006.
23. J. Niehren, L. Planque, J.-M. Talbot, and S. Tison. N-ary queries by tree automata. In *DBPL*, 2005.
24. N. Onose, A. Deutsch, Y. Papakonstantinou, and E. Curtmola. Rewriting nested XML queries using nested views. In *SIGMOD Conf.*, 2006.
25. T. Schwentick. On diving in trees. In *MFCS*, 2000.
26. T. Schwentick. Automata for XML - a survey. *J. Comput. Syst. Sci.*, 73(3):289–315, 2007.
27. B. ten Cate and L. Segoufin. XPath, transitive closure logic, and nested tree walking automata. In *PODS*, 2008.
28. J. W. Thatcher and J. B. Wright. Generalized finite automata theory with an application to a decision problem of second-order logic. *Mathematical Systems Theory*, 2(1):57–81, 1968.
29. A. Thomo and S. Venkatesh. Rewriting of VPLs for XML data integration. In *CIKM*, 2008.
30. J. Wang, J. Li, and J. X. Yu. Answering tree pattern queries using views: a revisit. In *EDBT*, 2011.
31. W. Xu and Z. M. Özsoyoglu. Rewriting XPath queries using materialized views. In *VLDB*, 2005.
32. C. Yu and L. Popa. Constraint-based XML query rewriting for data integration. In *SIGMOD*, 2004.

## 10 Appendix

### 10.1 Intersection

Given FTAs  $\mathcal{A}_1 = (S_1, \Sigma, F_1, \Delta_1)$  and  $\mathcal{A}_2 = (S_2, \Sigma, F_2, \Delta_2)$  for  $L_1$  and  $L_2$ , respectively, we construct FTA

$$\mathcal{A} = (S_1 \times S_2, \Sigma, F_1 \times F_2, \Delta_1 \otimes \Delta_2)$$

where  $\otimes$  is defined as follows.

We denote  $(s_1, s_2) \in S_1 \times S_2$  by  $s_{12}$ . Let  $H_1$  and  $H_2$  be two regular languages over  $S_1$  and  $S_2$ , respectively, We define

$$H_1 \otimes H_2 = \{s_{i_1 j_1} \dots s_{i_n j_n} : s_{i_1} \dots s_{i_n} \in H_1, \\ s_{j_1} \dots s_{j_n} \in H_2, \text{ and } n \in \mathbb{N}\}.$$

Also, we complete  $H_1 \otimes H_2$  by adding  $\epsilon$  to it, if  $\epsilon \in H_1$  and  $\epsilon \in H_2$ . Given  $\epsilon$ -free finite state automata (FSAs) for  $H_1$  and  $H_2$  one can easily build an FSA for  $H_1 \otimes H_2$  in polynomial time. Now, we define

$$\Delta_1 \otimes \Delta_2 = \{H_1 \otimes H_2 \xrightarrow{a} s_{12} : H_1 \xrightarrow{a} s_1 \text{ in } \Delta_1, \text{ and} \\ H_2 \xrightarrow{a} s_2 \text{ in } \Delta_2\}.$$

**Proposition 8.** (1)  $L(\mathcal{A}) = L(\mathcal{A}_1) \cap L(\mathcal{A}_2)$ . (2) Furthermore, (i) if  $\mathcal{A}_1$  and  $\mathcal{A}_2$  are complete, so is  $\mathcal{A}$ ; (ii) if  $\mathcal{A}_1$  and  $\mathcal{A}_2$  are deterministic, so is  $\mathcal{A}$ .

*Proof.* Let  $\rho$  be an accepting run of  $\mathcal{A}$  on a tree  $t$ . Notice  $\rho$  is a tree with the same shape as  $t$ , but with nodes labeled by states  $s_{ij} \in S_1 \times S_2$ . Consider now  $\rho_1$  obtained by changing each  $s_{ij}$  in  $\rho$  to  $s_i$ , and  $\rho_2$  obtained by changing each  $s_{ij}$  in  $\rho$  to  $s_j$ . It is easy to see  $\rho_1$  and  $\rho_2$  are accepting runs of  $\mathcal{A}_1$  and  $\mathcal{A}_2$  on  $t$ , respectively, and so  $L(\mathcal{A}) \subseteq L(\mathcal{A}_1) \cap L(\mathcal{A}_2)$ .

Conversely, let  $\rho_1$  and  $\rho_2$  be accepting runs of  $\mathcal{A}_1$  and  $\mathcal{A}_2$  on a tree  $t$ .  $\rho_1$  and  $\rho_2$  are trees with the same shape as  $t$ , but with nodes labeled by states in  $S_1$  and  $S_2$ , respectively. Now consider the tree  $\rho$  with the same shape as  $t$ , but with nodes labeled by elements  $s_{ij} \in S_1 \times S_2$ , such that  $s_i$  and  $s_j$  label the corresponding nodes in  $\rho_1$  and  $\rho_2$ , respectively. Clearly,  $\rho$  is an accepting run of  $\mathcal{A}$  on  $t$  and so,  $L(\mathcal{A}_1) \cap L(\mathcal{A}_2) \subseteq L(\mathcal{A})$ . This proves claim (1).

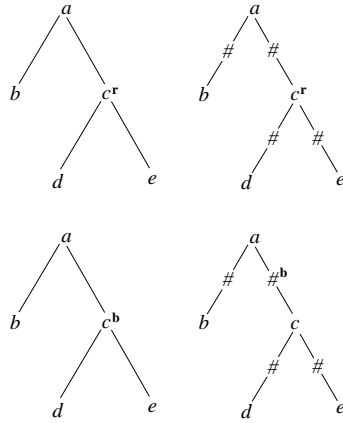
Furthermore, it is clear from the construction of  $\mathcal{A}$  and from that of runs  $\rho, \rho_1, \rho_2$  above that whenever  $\mathcal{A}_1$  and  $\mathcal{A}_2$  are complete, so is  $\mathcal{A}$ . Suppose now  $\mathcal{A}_1$  and  $\mathcal{A}_2$  are deterministic. From the construction, and since deterministic FTAs have at most one run on any tree data  $t$ , it follows that  $\mathcal{A}$  is deterministic as well.  $\square$

### 10.2 Avoiding the coincidence of query and view markings

**The transformation.** Let  $\# \notin \Sigma$ . If  $\mathbf{p} \in \Upsilon^r$ , we define  $\mathbf{p}\#$  to be the pattern obtained from  $\mathbf{p}$  by introducing a dummy node  $\xi_{xy}$  for each parent-child pair

$x, y$  of nodes in  $\mathfrak{p}$ , and making  $\xi_{xy}$  child of  $x$  and parent of  $y$ , and labeling  $\xi_{xy}$  by  $\#$ . For an example, see Figure 3, top. Furthermore, if  $\mathfrak{p}$  has the root red, then we create a new root labeled by  $\#$  with the old root as a child.

Let  $\#^b$  be the blue variant of  $\#$ . If  $\mathfrak{p} \in \mathcal{T}^b$ , we define  $\mathfrak{p}\#$  to be the pattern obtained from  $\mathfrak{p}$  by introducing as before a dummy node  $\xi_{xy}$ , for each parent-child pair  $x, y$  of nodes in  $\mathfrak{p}$ , and making  $\xi_{xy}$  child of  $x$  and parent of  $y$ . If  $y$  is black,  $\xi_{xy}$  is labeled by  $\#$ . Otherwise, i.e., if  $y$  is blue,  $\xi_{xy}$  is labeled by  $\#^b$ , and  $y$  changes color to black. For an example, see Figure 3, bottom. Furthermore, if  $\mathfrak{p}$  has the root blue, then we create a new root labeled by  $\#^b$  with the old root as a child. The old root is changed in this case to black.



**Fig. 3.** Transforming a pattern in  $\mathcal{T}^r$  (top), and a pattern in  $\mathcal{T}^b$  (bottom).

Let  $Q \subseteq \mathcal{T}^r$  and  $V \subseteq \mathcal{T}^b$ . We define

$$Q' = \{\mathfrak{p}\# : \mathfrak{p} \in Q\} \text{ and } V' = \{\mathfrak{p}\# : \mathfrak{p} \in V\}.$$

Clearly:

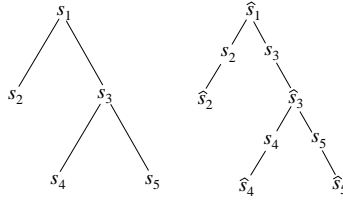
The red node in a  $Q'$  pattern will never coincide with the blue node in a matching  $V'$  pattern.

**Obtaining FTAs for  $Q'$  and  $V'$ .** Let  $\mathcal{A} = (S_{\mathcal{A}}, \Sigma \cup \Sigma^r, F_{\mathcal{A}}, \Delta_{\mathcal{A}})$  be a targeted FTA for  $Q$ , where  $S_{\mathcal{A}} = \{s_1, \dots, s_{n_{\mathcal{A}}}\}$ . Let  $\hat{S}_{\mathcal{A}} = \{\hat{s}_1, \dots, \hat{s}_{n_{\mathcal{A}}}\}$  be a set of new states, each corresponding to a state in  $S_{\mathcal{A}}$ .

We construct  $\mathcal{A}' = (S_{\mathcal{A}'}, \Sigma \cup \Sigma^r \cup \{\#\}, F_{\mathcal{A}'}, \Delta_{\mathcal{A}'})$ , where

$$\begin{aligned} S_{\mathcal{A}'} &= S_{\mathcal{A}} \cup \hat{S}_{\mathcal{A}} \\ F_{\mathcal{A}'} &= \{\hat{s}_i : s_i \in F_{\mathcal{A}}\} \cup \\ &\quad \{s_i : s_i \in F_{\mathcal{A}} \text{ and there exists } H \xrightarrow{a^r} s_i \text{ in } \Delta_{\mathcal{A}}\} \\ \Delta_{\mathcal{A}'} &= \{H \xrightarrow{a} \hat{s}_i \text{ and } \hat{s}_i \xrightarrow{\#} s_i : H \xrightarrow{a} s_i \text{ in } \Delta_{\mathcal{A}}\} \cup \\ &\quad \{H \xrightarrow{a^r} \hat{s}_i \text{ and } \hat{s}_i \xrightarrow{\#} s_i : H \xrightarrow{a^r} s_i \text{ in } \Delta_{\mathcal{A}}\}. \end{aligned}$$

We have that



**Fig. 4.** Corresponding runs of  $\mathcal{A}$  and  $\mathcal{A}'$ .

**Proposition 9.**  $\mathcal{A}'$  recognizes  $Q'$ .

*Proof.* The claim follows from the observation that there is a one-to-one correspondence between the runs of  $\mathcal{A}$  and the runs of  $\mathcal{A}'$ . An example is given in Figure 4. Let  $\rho, \varrho$  be corresponding runs. If  $\rho$  is a run of  $\mathcal{A}$  on a pattern  $p$ , then  $\varrho$  is a run  $\mathcal{A}'$  on  $p\#$ , and vice-versa.  $\square$

Let  $\mathcal{B} = (S_{\mathcal{B}}, \Sigma \cup \Sigma^b, F_{\mathcal{B}}, \Delta_{\mathcal{B}})$  be a targeted FTA for  $V$ , where  $S_{\mathcal{B}} = \{s_1, \dots, s_{n_{\mathcal{B}}}\}$ . Let  $\hat{S}_{\mathcal{B}} = \{\hat{s}_1, \dots, \hat{s}_{n_{\mathcal{B}}}\}$  be a set of new states, each corresponding to a state in  $S_{\mathcal{B}}$ . We construct  $\mathcal{B}' = (S_{\mathcal{B}'}, \Sigma \cup \Sigma^b \cup \{\#, \#^b\}, F_{\mathcal{B}'}, \Delta_{\mathcal{B}'})$ , where

$$\begin{aligned} S_{\mathcal{B}'} &= S_{\mathcal{B}} \cup \hat{S}_{\mathcal{B}} \\ F_{\mathcal{B}'} &= \{\hat{s}_i : s_i \in F_{\mathcal{B}}\} \cup \\ &\quad \{s_i : s_i \in F_{\mathcal{B}} \text{ and there exists } H \xrightarrow{a^b} s_i \text{ in } \Delta_{\mathcal{B}}\} \\ \Delta_{\mathcal{B}'} &= \{H \xrightarrow{a} \hat{s}_i \text{ and } \hat{s}_i \xrightarrow{\#} s_i : H \xrightarrow{a} s_i \text{ in } \Delta_{\mathcal{B}}\} \cup \\ &\quad \{H \xrightarrow{a} \hat{s}_i \text{ and } \hat{s}_i \xrightarrow{\#^b} s_i : H \xrightarrow{a^b} s_i \text{ in } \Delta_{\mathcal{B}}\}. \end{aligned}$$

We can show, similarly as for Proposition 9, that

**Proposition 10.**  $\mathcal{B}'$  recognizes  $V'$ .



**Rewriting.** Let  $R'$  be the MCR of  $Q'$  using  $V'$ . The automaton for  $R'$  should be transformed into an automaton over the original  $\Sigma \cup \Sigma^r$  in order to be useful to extract the answer from the view extension. For this, we will make use of “ $\epsilon$ ”-transition rules ( $\epsilon$ -TRs) of the form  $s' \rightarrow s$ .

Let  $\mathcal{C}' = (S_{\mathcal{C}'}, \Sigma \cup \Sigma^r \cup \{\#\}, F_{\mathcal{C}'}, \Delta_{\mathcal{C}'})$  be a targeted FTA for  $R'$ . By the definition of  $R'$ , each transition rule  $H \xrightarrow{\#} s$  of  $\mathcal{C}'$  should be such that  $H \subseteq S_{\mathcal{C}'}$ , i.e., the words of  $H$  are singleton states. Therefore, we can replace each  $H \xrightarrow{\#} s$  transition rule by a set of  $s' \xrightarrow{\#} s$  transition rules, for  $s' \in H$ .

Now we finally produce  $\mathcal{C}$  that is the same as  $\mathcal{C}'$ , but with the  $s' \xrightarrow{\#} s$  transition rules replaced by corresponding  $s' \rightarrow s$   $\epsilon$ -transition rules.  $\mathcal{C}$  is an FTA over  $\Sigma \cup \Sigma^r$ . It can be shown that FTAs with  $\epsilon$ -TRs can be converted into equivalent FTAs without  $\epsilon$ -TRs (cf. [10]).