

# The Complexity of Text-Preserving XML Transformations\*

Timos Antonopoulos  
Hasselt University and  
Transnational University of  
Limburg  
timos.antonopoulos@uhasselt.be

Wim Martens<sup>†</sup>  
TU Dortmund  
wim.martens@udo.edu

Frank Neven  
Hasselt University and  
Transnational University of  
Limburg  
frank.neven@uhasselt.be

## ABSTRACT

While XML is nowadays adopted as the de facto standard for data exchange, historically, its predecessor SGML was invented for describing electronic documents, i.e., marked-up text. Actually, today there are still large volumes of such XML texts. We consider simple transformations which can change the internal structure of documents, that is, the mark-up, and can filter out parts of the text but do not disrupt the ordering of the words. Specifically, we focus on XML transformations where the transformed document is a subsequence of the input document when ignoring mark-up. We call the latter *text-preserving* XML transformations. We characterize such transformations as copy- and rearrange-free transductions. Furthermore, we study the problem of deciding whether a given XML transducer is text-preserving over a given tree language. We consider top-down transducers as well as the abstraction of XSLT called DTL. We show that deciding whether a transformation is text-preserving over an unranked regular tree language is in PTIME for top-down transducers, EXPTIME-complete for DTL with XPath, and decidable for DTL with MSO patterns. Finally, we obtain that for every transducer in one of the above mentioned classes, the maximal subset of the input schema can be computed on which the transformation is text-preserving.

## Categories and Subject Descriptors

H.2.3 [Database Management]: Languages—*Data manipulation languages (DML)*; F.4.3 [Mathematical Logic and Formal Languages]: Formal Languages; D.2.4 [Software Engineering]: Software/Program Verification

\*We acknowledge the financial support of the Future and Emerging Technologies (FET) programme within the Seventh Framework Programme for Research of the European Commission, under the FET-Open grant agreement FOX, number FP7-ICT-233599.

<sup>†</sup>Supported by a grant of the North-Rhine Westfalian Academy of Sciences and Arts, and the Stiftung Mercator Essen.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

PODS'11, June 13–15, 2011, Athens, Greece.

Copyright 2011 ACM 978-1-4503-0660-7/11/06 ...\$10.00.

## General Terms

Algorithms, Theory, Verification

## 1. INTRODUCTION

While XML is embraced by the industry as the de facto standard for data exchange on the Web, XML can also be used to describe marked-up text. Actually, XML's predecessor, SGML, was exactly intended for this purpose. We refer to such documents as text-centric. Examples are poems, books, legislative text, e-governments text, and so on. A characteristic of such documents is that data values are words and sentences, and that their ordering matters. In this paper, we are interested in simple transformations which can change the internal structure of the document, that is, the mark-up, and can filter out parts of the text but do not disrupt the ordering of the words. In this way, the transformations can preserve the meaning of the filtered text. Specifically, we focus on XML transformations where the transformed document is a subsequence of the input document when ignoring mark-up. We call such transformations *text-preserving*. The central problem we consider in this paper, is to decide whether a given transformation is text-preserving for a given class of XML documents. Classes of XML documents will be defined by unranked regular tree languages and transformations by various kinds of transducers. We will show that, modulo some technical restrictions on transductions, transductions are text-preserving if and only if they are not copying and not rearranging.

We will not consider arbitrary transductions but will focus on XSLT-like transformations and use the abstraction of XSLT called DTL, introduced in [11], extended for dealing with data values. DTL is a rule-based language and is parameterized by a pattern language which is used for navigating in the tree and for deciding which rules can be executed. We consider three kinds of settings. The first is a top-down setting where navigation is restricted to children and rule patterns are restricted to label tests. Actually, to ease presentation, we define the top-down fragment of DTL separately as it corresponds to the formalism of top-down uniform tree transducers. In the two other settings, we consider XPath and MSO for both navigational and rule patterns. We refer to these fragments as  $DTL^{XPath}$  and  $DTL^{MSO}$ , respectively. We show that testing whether a transduction is text-preserving is in PTIME, in EXPTIME and decidable, respectively, for top-down transducers,  $DTL^{XPath}$  and  $DTL^{MSO}$  over the class of regular tree languages represented by non-deterministic tree automata.

The high-level proof idea underlying all three results is

the same, but the details differ greatly, leading to the different complexities. We essentially show that the set of trees for which the given transduction is not text-preserving is a regular tree language. We refer to the latter as the language of counter-examples. The result then follows by testing emptiness of the intersection of that language with the tree automaton representing the input schema. Of course, the language of counter-examples depends on the transducers. We represent it in the three respective settings by non-deterministic top-down tree automata, alternating tree-walking automata, and non-deterministic tree-jumping automata with MSO-transitions. Of independent interest, we observe that non-deterministic tree-jumping automata with MSO-transitions define only regular tree languages. As regular languages are closed under complement, it readily follows that a regular tree language can be constructed which represents the largest subset of the input schema on which a given transducer is text-preserving.

**Outline.** In Section 2, we introduce the necessary definitions. In Section 3, we characterize text-preserving transductions. In Section 4 and Section 5, we consider top-down uniform tree transducers and DTL, respectively. In Section 6, we discuss related work and we conclude in Section 7.

## 2. DEFINITIONS

For any  $n, m \in \mathbb{N}$  with  $m \geq n$ ,  $[n, m]$  denotes the set of integers  $\{n, n+1, \dots, m-1, m\}$ . In this paper we consider finite and infinite alphabets, but  $\Sigma$  always denotes a finite alphabet. We denote symbols from  $\Sigma$  by  $\sigma, \sigma_1$ , and so on. We assume that special symbols such as “(”, “)”, etc. are not in  $\Sigma$ .

**Strings.** A *symbol* is an element of a (finite or infinite) alphabet  $\Delta$  and a *string*  $w$  is a finite sequence of symbols  $\sigma_1 \cdots \sigma_n$  for some  $n \in \mathbb{N}$ . We define the *length* of a string  $w = \sigma_1 \cdots \sigma_n$ , denoted by  $|w|$ , to be  $n$  and we also refer to  $|w|$  as the *size* of  $w$ . The empty string is denoted by  $\varepsilon$ . We assume that readers are familiar with standard operations on strings and sets of strings such as concatenation.

**Automata.** A *nondeterministic finite string automaton (NFA)*  $A$  over an alphabet  $\Sigma$  is a tuple  $(Q, \Sigma, \delta, q_0, F)$ , such that  $Q$  is a finite set of states,  $q_0 \in Q$  is the initial state,  $F \subseteq Q$  is the set of final states and  $\delta : Q \times \Sigma \rightarrow 2^Q$  is the transition function. A *run* of the automaton  $A$  on a string  $w = a_1 \cdots a_n$  is a sequence of states  $q_0 q_1 \cdots q_n$  such that, for each  $i \in [1, n]$ ,  $q_i \in \delta(q_{i-1}, a_i)$ . A run on  $w$  is *accepting* if  $q_n \in F$ , and a string  $w$  is *accepted* by  $A$ , if there exists an accepting run of  $A$  on  $w$ . The *size*  $|A|$  of an NFA  $A$  is its total number of states and transitions. By  $\mathcal{L}(A)$  we denote the set of strings accepted by  $A$ . A string language  $L$  over  $\Sigma$  is *regular* if there is an NFA  $A$  such that  $\mathcal{L}(A) = L$ . We denote the set of regular languages over alphabet  $\Sigma$  by  $\text{REG}(\Sigma)$ .

**Hedges and trees.** The set of *unranked trees* over alphabet  $\Sigma$ , denoted by  $\text{Trees}_\Sigma$ , is the smallest set  $S$  of strings over  $\Sigma$  and the parenthesis symbols “(” and “)” such that,  $\varepsilon \in S$  and for each  $\sigma \in \Sigma$  and  $w \in S^*$ , we have that  $\sigma(w)$  is in  $S$ . For readability, we denote the tree  $\sigma()$  by  $\sigma$ . A *tree language* over  $\Sigma$  is a subset of  $\text{Trees}_\Sigma$ . The set of *unranked hedges* over alphabet  $\Sigma$ , denoted by  $\text{Hedges}_\Sigma$ , is defined as  $\text{Hedges}_\Sigma = \text{Trees}_\Sigma^*$ . In particular, each tree is also a hedge.

When we write a tree  $t$  as  $t = \sigma(t_1 \cdots t_n)$  or a hedge  $h$  as  $h = t_1 \cdots t_n$ , we implicitly assume that all  $t_i$  are trees.

The *set of nodes* of a tree  $t$  and of a hedge  $h$ , denoted by  $\text{Nodes}^t$  and  $\text{Nodes}^h$ , respectively, are sets of strings in  $\mathbb{N}^*$  and are inductively defined as follows. If  $h = \sigma_1 \cdots \sigma_n$ , then  $\text{Nodes}^h = \{1, \dots, n\}$ . Here, for each  $i \in [1, n]$ , the node  $i$  is labelled by  $\sigma_i$ , which we denote by  $\text{lab}^h(i) = \sigma_i$ . If  $t = \sigma(h)$ , then  $\text{Nodes}^t = \{1\} \cup \{1u \mid u \in \text{Nodes}^h\}$ . Here, the root node is 1 and the other nodes are nodes of the subhedge  $h$ , prefixed by a 1. The root is labeled  $\sigma$ , i.e.,  $\text{lab}^t(1) = \sigma$  and, for every node  $1u \in \text{Nodes}^t$  we define  $\text{lab}^t(1u) = \text{lab}^h(u)$ . Finally, if  $h = t_1 \cdots t_n$ , then  $\text{Nodes}^h = \cup_{i=1}^n \{iu \mid u \in \text{Nodes}^{t_i}\}$ . Here, the labels from the  $t_i$  carry over to  $h$ , that is, for each  $i \in [1, n]$  and each node  $iu \in \text{Nodes}^h$ , we have  $\text{lab}^h(iu) = \text{lab}^{t_i}(iu)$ . The *depth* of a node  $u \in \mathbb{N}^*$  is  $|u|$ . Hence, the depth of the root of a tree is one.

From the definition of nodes we can see that the lexicographic order  $<_{\text{lex}}$  on  $\text{Nodes}^h$  corresponds to the order generated by the depth-first (left-to-right) traversal of the hedge  $h$ . More specifically, for two nodes  $iu$  and  $ju \in \text{Nodes}^h$  with  $i, j \in \mathbb{N}$ , we have  $iu <_{\text{lex}} ju$  if  $i < j$  or if  $i = j$  and  $u <_{\text{lex}} v$ . The *children* of a node  $v$  in  $h$  are all nodes  $v' \in \text{Nodes}^h$  such that  $v' = v \cdot i$  for  $i \in \mathbb{N}$ . A node  $v \in \text{Nodes}^h$  is a *leaf* if it has no children.

The *size* of a hedge  $h$ , denoted by  $|h|$ , is its number of nodes. For a tree  $t$  and a node  $u \in \text{Nodes}^t$ , we denote by  $\text{anc-str}^t(u)$  the *ancestor string* of  $u$  in  $t$ , i.e., the string formed by the labels on the path in the tree  $t$  from the root to  $u$ , including the label of  $u$ . The *lowest common ancestor* of two nodes  $v_1$  and  $v_2$  in  $\text{Nodes}^t$  is the node corresponding to the longest common prefix of  $v_1$  and  $v_2$ .

The *subtree of hedge  $h$  at a node  $u$*  is the tree induced by the set of nodes with prefix  $u$  and is denoted by  $\text{subtree}^h(u)$ . For any hedge  $h$ , node  $u$ , and for any hedge  $h'$ ,  $h[u \leftarrow h']$ , denotes the hedge obtained from  $h$  by replacing  $\text{subtree}^h(u)$  with  $h'$  and by redefining the set of nodes and the label function accordingly. For any two alphabets  $\Sigma$  and  $\Gamma$ ,  $\text{Trees}_\Sigma(\Gamma)$  is the set of trees over the alphabet  $\Sigma \cup \Gamma$ , where only leaves are allowed to be labelled with symbols from  $\Gamma$ . Similarly,  $\text{Hedges}_\Sigma(\Gamma)$  is  $(\text{Trees}_\Sigma(\Gamma))^*$ . Finally, the *frontier*<sup>1</sup> of a hedge  $h$ , denoted by  $\text{frontier}(h)$ , is the largest sequence  $\text{lab}^h(v_1) \cdots \text{lab}^h(v_n)$ , where, for every  $i \in [1, n]$ ,  $v_i$  is a leaf of  $h$ , and for every  $i \in [1, n-1]$ ,  $v_i <_{\text{lex}} v_{i+1}$ .

In the following, we will often write  $\text{Nodes}$ ,  $\text{lab}$ , etc. without index whenever the tree or hedge is clear from the context.

**Text trees.** Let  $\text{Text}$  be an infinite set, disjoint from any finite alphabet we consider here, such as  $\Sigma$ . A *text tree*  $t$  over  $\Sigma$  is a tree in  $\text{Trees}_\Sigma(\text{Text})$  or, in other words, a tree over the alphabet  $\Sigma \cup \text{Text}$  where symbols from  $\text{Text}$  can only appear at the leaves. We refer to these nodes as *text nodes* and denote the set of text nodes of a tree  $t$  by  $\text{text-nodes}^t$ . A text tree language is a subset of  $\text{Trees}_\Sigma(\text{Text})$ . In the following, unless otherwise stated, whenever we say *tree* we always mean *text tree*. In particular, we also simply say *tree language* instead of *text tree language*.

A tree language  $L$  is *closed under Text-substitutions* if, for any tree  $t \in L$  and any text node  $u$  of  $t$ , the tree obtained from  $t$  by changing  $u$ 's label to another value in  $\text{Text}$  is also in the language  $L$ . Formally, a *Text-substitution*  $\rho$  is a function from  $\text{Trees}_\Sigma(\text{Text})$  to  $\text{Trees}_\Sigma(\text{Text})$ , such that

<sup>1</sup>A frontier is sometimes also called a *yield*.

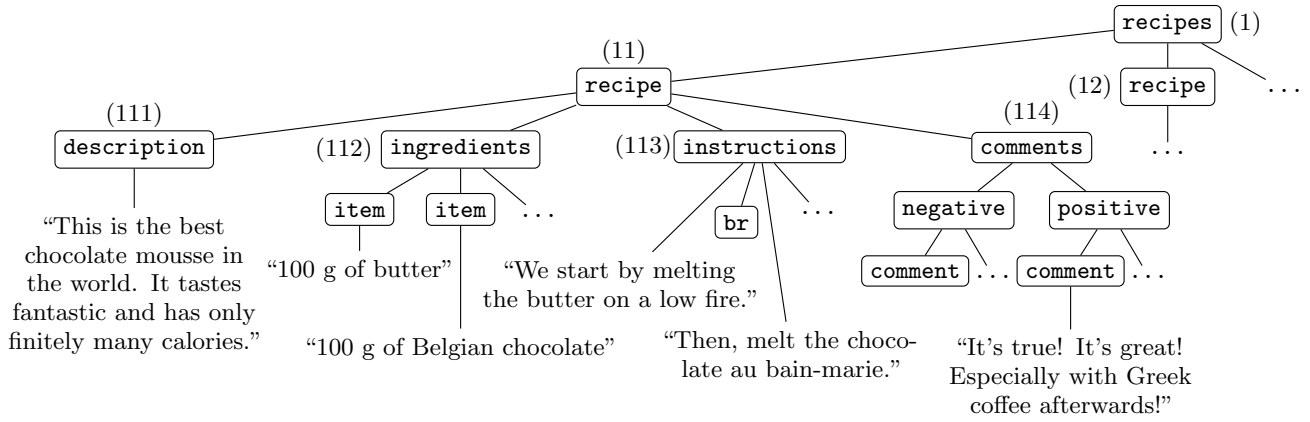


Figure 1: A text tree representing an XML document for recipes.

for each  $t \in \text{Trees}_\Sigma(\text{Text})$ ,  $\text{Nodes}^t = \text{Nodes}^{\rho(t)}$ , for each node  $v \notin \text{text-nodes}^t$ ,  $\text{lab}^t(v) = \text{lab}^{\rho(t)}(v)$  and for each node  $v \in \text{text-nodes}^t$ ,  $\text{lab}^{\rho(t)}(v) \in \text{Text}$ . In other words, a **Text**-substitution relabels zero or more leaf nodes labelled with a label in **Text**, to some other label in **Text**.

In the sequel we only consider tree languages that are closed under **Text**-substitutions primarily because we want to consider **Text** as an abstract set. Furthermore, when ignoring integrity constraints, XML schema languages are rather restricted in enforcing constraints on actual text values. Therefore, it makes sense to simply consider one infinite set of possible values so that replacing a certain text value with another does not change validity of the document with respect to the schema.

The *text-content* of a tree  $t$ , denoted by  $\text{text-content}(t)$ , is the string over alphabet **Text** obtained by concatenating the text values of all text nodes in document order, which is the lexicographic order induced by  $\mathbb{N}^*$ . Notice that, for any tree  $t$ ,  $\text{text-content}(t)$  is the substring of  $\text{frontier}(t)$  containing exactly the labels from **Text**.

EXAMPLE 2.1. Figure 1 depicts a tree representation of an XML document underlying a web site with recipes. Each recipe has a description, a list of ingredients, instructions, and a list of comments by users. The text nodes of the tree are the ones that contain text in quotation marks. The ancestor path of the node labeled **positive** is **recipes recipe comments positive**. The text-content of the tree is the concatenation of all the text between the quotation marks, from left to right. For clarity, we annotate some nodes with their name in braces.

A *tree transduction*  $T$  is a mapping  $T : L_1 \rightarrow L_2$  for tree languages  $L_1$  and  $L_2$ . We say that a string  $s_1 = \sigma_1 \cdots \sigma_n$  over alphabet  $\Delta$  is a *subsequence* of  $s_2$ , denoted  $s_1 \prec s_2$ , if  $s_2$  is of the form  $w_0 \sigma_1 w_1 \cdots w_{n-1} \sigma_n w_n$  for some  $w_0, \dots, w_n \in \Delta^*$ . We are now ready to define the notion central to this paper.

DEFINITION 2.2. A tree transduction  $T$  is *text preserving* over a tree language  $L$  if, for all trees  $t \in L$ ,

$$\text{text-content}(T(t)) \prec \text{text-content}(t).$$

**Schema languages.** We abstract the schema languages *Document Type Definition (DTD)* and *Relax NG* by ex-

tended context-free grammars and by unranked tree automata, respectively.

A *Document Type Definition (DTD)* over some finite alphabet  $\Sigma$  is a tuple  $D = (\Sigma \uplus \{\mathbf{text}\}, C, d, S_d)$  where  $C$  is a set of regular string languages over  $\Sigma \uplus \{\mathbf{text}\}$ ,  $d$  is a function that maps symbols in  $\Sigma$  to languages in  $C$ , and  $S_d \subseteq \Sigma$  is a set of start symbols. We refer to the languages in  $C$  as the *content models* of the DTD. A tree  $t$  is *valid with respect to* a DTD  $D$  or *satisfies*  $D$ , if its root is labelled by an element of  $S_d$  and, for every node labelled with some  $\sigma \in \Sigma$ , the sequence  $\sigma_1 \cdots \sigma_n$  of labels of its children, where any element in **Text** is replaced by the symbol **text**, is in the language  $d(\sigma)$ . So, **text** serves as a placeholder for text nodes. In the remainder of this paper, we assume that the regular languages in  $C$  are represented by regular expressions or NFAs.

The set of trees that are valid with respect to a DTD  $D$  is denoted by  $\mathcal{L}(D)$ . A DTD  $D$  is *reduced* if, for every  $\sigma \in \Sigma$  for which  $d(\sigma)$  is defined, there exists a tree  $t \in \mathcal{L}(D)$  such that the label  $\sigma$  occurs somewhere in  $t$ . Any DTD can be transformed to an equivalent reduced DTD in polynomial time [1, 16]. However, reducing a DTD is PTIME-complete.<sup>2</sup> In the following, we assume that all DTDs are reduced.

EXAMPLE 2.3. The tree in Figure 1 is valid w.r.t. the DTD  $(\Sigma \uplus \{\mathbf{text}\}, C, d, S_d)$ , where  $\Sigma$  is the set of labels used in the boxed nodes of the tree, i.e., **recipes**, **recipe**, ... We represent the regular languages in  $C$  by regular expressions. Then,  $S_d = \{\mathbf{recipes}\}$  and  $d$  is defined as:

<b>recipes</b>	$\mapsto$	<b>recipe</b> *
<b>recipe</b>	$\mapsto$	<b>description</b> · <b>ingredients</b> · <b>instructions</b> · <b>comments</b>
<b>ingredients</b>	$\mapsto$	<b>item</b> *
<b>instructions</b>	$\mapsto$	<b>(br + text)</b> *
<b>br</b>	$\mapsto$	$\varepsilon$
<b>comments</b>	$\mapsto$	<b>negative</b> · <b>positive</b>
<b>positive</b>	$\mapsto$	<b>comment</b> *
<b>negative</b>	$\mapsto$	<b>comment</b> *
<b>description</b>	$\mapsto$	<b>text</b>
<b>item</b>	$\mapsto$	<b>text</b>
...		

Notice that the children of **instructions** can come from  $\Sigma$

<sup>2</sup>More accurately, *deciding* whether a DTD is reduced is PTIME-complete.

and from **Text**. For every  $\Sigma$ -symbol  $\sigma$  for which we did not yet define  $d$ , we define  $d(\sigma) = \mathbf{text}$ .

Relax NG schemas are abstracted by *nondeterministic unranked tree automata* (NTA). An NTA  $N$  over some finite alphabet  $\Sigma$ , is a tuple  $(Q, \Sigma \uplus \{\mathbf{text}\}, \delta, q_0, F)$  where  $Q$  is a finite set of states,  $q_0 \in Q$  is the initial state,  $F$  is the set of final states, and  $\delta : Q \times (\Sigma \uplus \{\mathbf{text}\}) \rightarrow \text{REG}(Q)$  is the transition function. A *run* of  $N$  over some tree  $t$  is a function  $\rho : \text{Nodes}^t \rightarrow Q$  such that  $\rho(1) = q_0$ , i.e., the root of  $t$  is assigned  $q_0$  and, for any node  $v$  labeled with  $\sigma \in \Sigma$  with  $n$  children, it holds that  $\rho(v1) \cdots \rho(vn)$  is in the language  $\delta(\rho(v), \sigma)$ . A run  $\rho$  of  $N$  on some tree  $t$  is *accepting* if  $\varepsilon \in \delta(\rho(v), \mathbf{lab}(v))$  for every leaf  $v$  of  $t$  if  $\mathbf{lab}(v) \in \Sigma$ , and  $\delta(\rho(v), \mathbf{text}) = \{\varepsilon\}$  if  $\mathbf{lab}(v) \in \mathbf{Text}$ . So, also here **text** serves as a placeholder for **Text**-values. We define  $F$  to be the set of states  $q$  such that  $\varepsilon \in \delta(q, a)$  for some  $a \in \Sigma \uplus \{\mathbf{text}\}$ . A tree  $t$  is *accepted* by  $N$  if there exists an accepting run of  $N$  on  $t$ . By  $\mathcal{L}(N)$  we denote the set of trees accepted by  $N$ . A tree language  $L$  is *regular*, if there is an NTA  $N$  such that  $L = \mathcal{L}(N)$ . It is well-known that regular tree languages are closed under union, intersection, and complementation.

Unless otherwise mentioned, we assume that the regular languages  $\delta(q, \sigma)$  in NTAs are represented by NFAs. The size  $|N|$  of an NTA  $N = (Q, \Sigma, \delta, q_0, F)$  is equal to  $|Q| + |\delta|$ , where  $|\delta| = \sum_{q \in Q, a \in \Sigma} |\delta(q, a)|$ .

**Central problem.** In the rest of this paper, we consider the following decision problem. Let  $\mathcal{T}$  be a class of tree transductions and let  $\mathcal{L}$  be a class of tree languages. We then study the question:

Given a language  $L \in \mathcal{L}$  and a transduction  $T \in \mathcal{T}$ , is  $T$  text-preserving over  $L$ ?

Recall that, in this paper, the languages  $\mathcal{L}$  are always closed under **Text**-substitutions. In the conclusions we discuss how our technique also allows to solve more complex decision problems.

### 3. A CHARACTERIZATION OF TEXT-PRESERVING TRANSDUCTIONS

We provide a characterization of when a transduction is text-preserving in terms of its copying and rearranging behavior. First, we need some terminology. We say that a tree is *value-unique* when all its **Text**-values occurring at leaves are different. Notice that, since we only consider tree languages that are closed under **Text**-substitutions, all the tree languages we consider contain at least one tree which is value-unique.

**DEFINITION 3.1.** A transduction  $T$  is *copying over a tree language*  $L$  if there is a value-unique  $t \in L$  such that  $T(t)$  contains multiple occurrences of the same **Text**-value. A transduction  $T$  is *rearranging over a tree language*  $L$  if there is a value-unique  $t \in L$  such that, for some  $\gamma_1$  and  $\gamma_2 \in \mathbf{Text}$ ,  $\gamma_1\gamma_2 \prec \mathbf{text-content}(t)$  and  $\gamma_2\gamma_1 \prec \mathbf{text-content}(T(t))$ .

Up to now, transductions are general mappings from trees to trees. We will next put some restrictions on them. Let  $\mathfrak{X} \notin \mathbf{Text}$ . For  $\gamma \in \mathbf{Text} \cup \{\mathfrak{X}\}$ , let  $\rho_\gamma$  be the **Text**-substitution that for every tree  $t$ , relabels every text node with  $\gamma$ . A transduction  $T$  is *Text-independent* if for every tree  $t$  and every **Text**-substitution  $\rho$  on  $t$ ,

$$\rho_{\mathfrak{X}}(T(\rho(t))) = \rho_{\mathfrak{X}}(T(t)).$$

Informally, this ensures that the structure of the transduced tree  $T(t)$  does not depend on the **Text**-values occurring in  $t$ . Only the concrete **Text**-values in  $T(t)$  and  $T(\rho(t))$  can differ. The final notion we need is that of a **Text**-functional transduction. A transduction  $T$  is *Text-functional* if, for every tree  $t$ , there exists a function  $f$  from the set of text nodes of  $T(t)$  to the set of text nodes of  $t$ , such that for every **Text**-substitution  $\rho$  on  $t$ ,  $\mathbf{lab}^{T(\rho(t))}(v) = \mathbf{lab}^{\rho(t)}(f(v))$  for every node  $v \in \text{text-nodes}^{T(\rho(t))}$ . Intuitively, this means that **Text**-values in the output tree are determined by the **Text**-value at the corresponding node, determined by  $f$ , in the input tree. Note that a **Text**-functional transduction can never introduce **Text**-values that do not appear in the input tree.

**DEFINITION 3.2.** A transduction is *admissible* if it is **Text**-independent and **Text**-functional.

We will show in the sequel that all transductions we consider in this paper are admissible. We are now ready to prove the characterization of text-preserving transductions.

**THEOREM 3.3.** *An admissible transduction  $T$  is text-preserving over a language  $L$  if and only if it is not copying and not rearranging over  $L$ .*

**PROOF.** When  $T$  is copying or rearranging over  $L$  then  $T$  is obviously not text-preserving over  $L$ .

When  $T$  is not text-preserving, there exists a tree  $t \in L$  such that  $\mathbf{text-content}(T(t)) \not\prec \mathbf{text-content}(t)$ . Let  $\rho$  be a **Text**-substitution such that  $\rho(t)$  is value-unique. First, we argue that  $\mathbf{text-content}(T(\rho(t))) \not\prec \mathbf{text-content}(\rho(t))$ . Assume for the sake of contradiction that  $\mathbf{text-content}(T(\rho(t))) \prec \mathbf{text-content}(\rho(t))$  and let  $g$  be a function from  $\text{text-nodes}^{T(\rho(t))}$  to  $\text{text-nodes}^{\rho(t)}$  that witnesses this subsequence relation. In other words, since  $\rho(t)$  is value-unique,  $g$  maps every text node in  $T(\rho(t))$  to the unique text node in  $\rho(t)$  with the same text value. In particular,  $g$  is unique.

Since  $T$  is **Text**-functional by assumption, there exists a function  $f : \text{text-nodes}^{T(t)} \rightarrow \text{text-nodes}^t$  such that, for every **Text**-substitution  $\rho'$  and every node  $v \in \text{text-nodes}^{T(\rho'(t))}$ ,  $\mathbf{lab}^{T(\rho'(t))}(v) = \mathbf{lab}^{\rho'(t)}(f(v))$ . Since this condition holds for every substitution  $\rho'$ , it holds for  $\rho$  in particular. Therefore it is also the case that, for every node  $v \in \text{text-nodes}^{T(\rho(t))}$ ,  $\mathbf{lab}^{T(\rho(t))}(v) = \mathbf{lab}^{\rho(t)}(f(v))$ . Since  $\rho(t)$  is value unique and  $T$  is **Text**-independent, it follows that  $f$  maps every text node in  $T(\rho(t))$  to the unique text node in  $\rho(t)$  with the same text value.

Since  $\text{text-nodes}^t = \text{text-nodes}^{\rho(t)}$  and  $\text{text-nodes}^{T(t)} = \text{text-nodes}^{T(\rho(t))}$  because  $T$  is **Text**-independent, it follows that  $f$  and  $g$  are the same function. By definition of  $f$ , for every node  $v \in \text{text-nodes}^{T(t)}$ ,  $\mathbf{lab}^{T(t)}(v) = \mathbf{lab}^t(f(v))$  and therefore,  $f$  is a witness function to  $\mathbf{text-content}(T(t)) \prec \mathbf{text-content}(t)$ , which contradicts our assumption that  $T$  is not text-preserving on  $t$ .

It remains to show that  $T$  is either copying or rearranging on  $t$ . Consider the function  $f$ , and notice that since  $\mathbf{text-content}(T(t)) \not\prec \mathbf{text-content}(t)$ ,  $f$  is either not injective or does not preserve the order  $\prec_{\text{lex}}$  in trees. In the first case,  $T$  is copying and in the second case  $T$  is rearranging.  $\square$

### 4. TOP-DOWN TRANSDUCERS

We start with the simple uniform top-down tree transducers considered in [13, 14, 15] in the context of type checking.

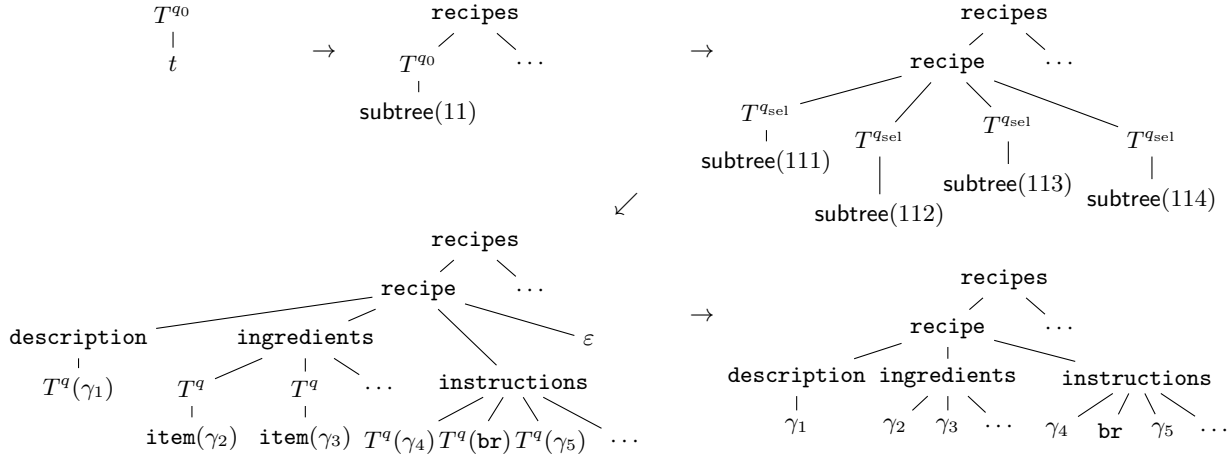


Figure 2: A uniform transducer selecting all descriptions, ingredients, and instructions from the tree of Fig. 1.

These correspond to a simple top-down fragment of XSLT and are equivalent to DTL, defined in the next section, instantiated with only downward navigation.

#### 4.1 Definition

DEFINITION 4.1. A *top-down uniform tree transducer* is a tuple  $T = (Q, \Sigma \cup \{\text{text}\}, q^0, R)$ , where  $Q$  is a finite set of states,  $\Sigma$  is a finite alphabet,  $q^0 \in Q$  is the initial state, and  $R$  is a finite set of rules of the form  $(q, a) \rightarrow h$ , where  $q \in Q$  and

- (a) if  $a \in \Sigma$  then  $h \in \text{Hedges}_\Sigma(Q)$  and
- (b) if  $a = \text{text}$  then  $h$  is **text**.

When  $q = q^0$ , then  $h$  is restricted to  $\text{Trees}_\Sigma(Q)$  and should contain at least one  $\Sigma$ -label. This restriction of  $h$  ensures that the output of  $T$  is always a tree. For each state  $q \in Q$  and symbol  $a \in \Sigma$  there is at most one hedge  $h \in \text{Hedges}_\Sigma(Q)$  such that  $(q, a) \rightarrow h$  is a rule in  $R$ .

The translation defined by  $T$  on a tree  $t$  in state  $q$ , denoted by  $T^q(t)$ , is defined inductively as follows:

- (i)  $T^q(\varepsilon) = \varepsilon$ ;
- (ii) for each  $\gamma \in \text{Text}$ , if there is a rule  $(q, \text{text}) \rightarrow \text{text}$  then  $T^q(\gamma) = \gamma$ , if there is no such rule, then  $T^q(\gamma) = \varepsilon$ ;
- (iii) if  $t = a(t_1 \cdots t_n)$  for some  $a \in \Sigma$  and there is a rule  $(q, a) \rightarrow h$ , then  $T^q(t)$  is obtained from  $h$  by replacing every node  $u$  in  $h$  labelled with  $p \in Q$ , by the hedge  $T^p(t_1) \cdots T^p(t_n)$ . If there is no such rule, then  $T^q(t) = \varepsilon$ .

The transformation of  $t$  by  $T$ , denoted by  $T(t)$ , is defined as  $T^{q^0}(t)$ . To make use of our characterization in Theorem 3.3, we need to be able to reason about the function that is defined by  $T$ . To this end, we say that the *transduction defined by  $T$*  or *transduction of  $T$*  is the function mapping every tree  $t$  onto  $T(t)$ . For any rule  $(q, a) \rightarrow h$  in  $R$ , we denote the hedge  $h$  on its right hand side by  $\text{rhs}(q, a)$ . The *size* of  $T$ , denoted by  $|T|$ , is equal to  $|Q| + |R|$ , where  $|R| = \sum_{q \in Q, a \in \Sigma} |\text{rhs}(q, a)|$ . In the remainder of this section, we simply say transducer rather than top-down uniform tree transducer.

EXAMPLE 4.2. The uniform transducer defined by the rules

$$\begin{aligned} (q_0, \text{recipes}) &\rightarrow \text{recipes}(q_0) \\ (q_0, \text{recipe}) &\rightarrow \text{recipe}(q_{\text{sel}}) \\ (q_{\text{sel}}, \sigma) &\rightarrow \sigma(q) \quad (\sigma \in \{\text{description, ingredients, instructions}\}) \end{aligned}$$

$$\begin{aligned} (q, \text{item}) &\rightarrow q \\ (q, \text{br}) &\rightarrow \text{br}(q) \\ (q, \text{text}) &\rightarrow \text{text} \end{aligned}$$

Selects, from the tree of Figure 1, all recipes, their descriptions, ingredients lists, and instructions; and deletes the comments. Furthermore, it keeps the markup (i.e., the **br** nodes) in the instructions, but deletes the **item** nodes. The transformation on the tree in Figure 1 is illustrated in Figure 2. For readability in the figure, we sometimes write a function call of the form  $T^q(t)$  as a (sub)tree where  $t$  is a child of a node labelled  $T^q$ .

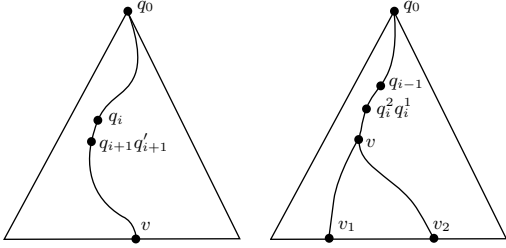
We call a state  $q$  of  $T$  *reachable* if there is a sequence of pairs  $(q_1, a_1) \cdots (q_n, a_n)$  for some  $n \in \mathbb{N}$ , such that  $q_1 = q^0$ ,  $q_n = q$  and, for all  $i \in [1, n - 1]$ , we have that  $q_{i+1}$  occurs as the label of a leaf in the hedge  $\text{rhs}(q_i, a_i)$ . Notice that, if a transducer has a rule of the form  $(q, a) \rightarrow \varepsilon$ , we can remove this rule and the resulting transducer is equivalent (due to rule (iii) in the definition of  $T^q(t)$ ). We therefore say that a rule of the form  $(q, a) \rightarrow \varepsilon$  is *useless*. A transducer is *reduced*, if all its states are reachable and it contains no useless rules. In the following, we assume that transducers are always reduced.

Notice that rules in transducers do not contain any values from **Text**. This ensures that transducers can not introduce a **Text**-value not present in the input tree. We note that the transducers in this section are admissible, in the sense of Definition 3.2.

LEMMA 4.3. *Top-down uniform tree transducers are admissible.*

#### 4.2 Copying and rearranging transducers

In this section, we provide equivalent formalizations for the notions of copying and rearranging introduced in Section 3. In Section 4.3, we then show that it is decidable in PTIME whether a transducer is copying or rearranging thereby obtaining PTIME-decidability for testing whether a transducer is text-preserving.



**Figure 3: Illustration of copying (left) and rearranging (right) for top-down uniform tree transducers.**

First we formally define copying and rearranging for a transducer.

**DEFINITION 4.4.** Let  $L$  be a tree language. A transducer  $T$  is *copying over  $L$*  if the transduction defined by  $T$  is copying over  $L$ . It is *rearranging over  $L$*  if the transduction defined by  $T$  is rearranging over  $L$ .

In the following lemma, we provide an operational condition that is equivalent to the general notion of copying defined above. First, we introduce the following definitions. A *text path* is a sequence in  $\Sigma^* \cdot \text{Text}$ . A *path run* of  $T$  on a text path  $a_1 \cdots a_n \gamma$  is a sequence  $q_1, \dots, q_n, q_{n+1}$  such that  $q_1 = q^0$ , and for all  $i \in [1, n]$ ,  $q_{i+1}$  occurs at a leaf of  $\text{rhs}(q_i, a_i)$  and  $(q_{n+1}, \text{text}) \rightarrow \text{text}$ . Intuitively, a path run is a set of states that  $T$  assumes when processing the text path. Note that there can be several path runs for the same text path as the transducer can use rules containing multiple states in a right-hand side.

**LEMMA 4.5.** A transducer  $T = (Q, \Sigma \cup \{\text{text}\}, q^0, R)$  is copying over a tree language  $L$  if and only if there exists a tree  $t \in L$  and a leaf node  $v$  in  $t$  such that  $\text{anc-str}(v) = a_1 \cdots a_n \cdot \gamma$  is a text path on which either

- (1)  $T$  has two different path runs, or
- (2)  $T$  has a path run  $q_1, \dots, q_n, q_{n+1}$ , where for some  $i \in [1, n]$ ,  $q_{i+1} \cdot q_{i+1} \prec \text{frontier}(\text{rhs}(q_i, a_i))$ .

Recall that  $\prec$  denotes the subsequence relation. Condition (2) therefore states that  $q_{i+1}$  occurs on two different leaves of  $\text{rhs}(q_i, a_i)$ . The condition in Lemma 4.5 is illustrated in Figure 3. In the figure, if  $q_{i+1} \neq q'_{i+1}$ , there are two path runs and if  $q_{i+1} = q'_{i+1}$ , then  $q_{i+1} \cdot q_{i+1} \prec \text{frontier}(\text{rhs}(q_i, a_i))$ .

Next, we characterize when a transducer rearranges two Text-values in a tree  $t$ . Similarly as for copying, we provide an operational condition for rearranging. Intuitively, a transducer rearranges when there are two leaves  $v_1$  and  $v_2$  in  $t$  and  $T(t)$  swaps the Text-values corresponding to  $v_1$  and  $v_2$ . In terms of a top-down uniform tree transducer, this means that, on the path from the root to the lowest common ancestor of  $v_1$  and  $v_2$ , a swap operation takes place.

**LEMMA 4.6.** A transducer  $T = (Q, \Sigma \cup \{\text{text}\}, q^0, R)$  is rearranging over a tree language  $L$  if and only if there exists a tree  $t \in L$ , two leaf nodes  $v_1, v_2$  of  $t$  with  $v_1 \prec_{\text{lex}} v_2$  with text paths  $\text{anc-str}(v_1) = a_1^1 \cdots a_n^1 \cdot \gamma^1$  and  $\text{anc-str}(v_2) = a_1^2 \cdots a_m^2 \cdot \gamma^2$  on which which  $T$  has path runs  $q_1^1, \dots, q_n^1, q_{n+1}^1$  and  $q_1^2, \dots, q_m^2, q_{m+1}^2$ , such that the following holds. If the lowest common ancestor  $v$  of  $v_1$  and  $v_2$  has depth  $k$ , then there exists an  $i \leq k$  such that

- (1) for all  $j < i$ , we have  $(q_j^1, a_j^1) = (q_j^2, a_j^2)$  and
- (2)  $q_i^2 \cdot q_i^1 \prec \text{frontier}(\text{rhs}(q_{i-1}^1, a_{i-1}^1))$ .

The condition in Lemma 4.6 for rearranging is illustrated in Figure 3. From Lemmas 4.5 and 4.6 and from Theorem 3.3 we now obtain the following characterization for text-preserving top-down uniform transducers:

**THEOREM 4.7.** A uniform top-down transducer  $T$  is not text-preserving over a tree language  $L$  if and only if the condition in Lemma 4.5 or the condition in Lemma 4.6 holds.

### 4.3 PTIME Result

This section is devoted to the proof of Theorem 4.11 which states that it is decidable in PTIME whether a top-down tree transducer is text-preserving over an unranked regular tree language represented by an NTA.

For a tree language  $L$ , we denote by  $L_{\text{text}}$  the language over alphabet  $\Sigma \uplus \{\text{text}\}$  obtained from  $L$  by replacing, in each tree, each label  $\gamma \in \text{Text}$  by the label  $\text{text}$ . The *path language* of a regular tree language  $L$  is the language of all root-to-leaf paths in trees of  $L_{\text{text}}$  that end with  $\text{text}$ . More formally, the text path language of  $L$  is  $\{w \mid t \in L_{\text{text}}, \text{lab}^t(u) = \text{text}, \text{ and } w = \text{anc-str}^t(u)\}$ . We say that an NFA  $A_L$  is a *path automaton for regular tree language  $L$*  if  $A_L$  accepts the path language of  $L$ . For a tree transducer  $T$ , we say that an NFA  $A_T$  is a *transducer path automaton for  $T$*  if it accepts exactly the text paths on which  $T$  has a path run. Notice that  $A_L$  and  $A_T$  accept only strings that end with  $\text{text}$ . We show that they can be constructed in PTIME.

- LEMMA 4.8.** (1) Given an NTA  $N$ , we can construct a path automaton  $A_N$  for  $\mathcal{L}(N)$  in polynomial time.  
(2) Given a uniform tree transducer  $T$ , we can construct a transducer path automaton  $A_T$  for  $T$  in polynomial time.

From now on, we refer to  $A_T$  as the transducer path automaton of  $T$  and to  $A_N$  as the path automaton of  $N$ . As our next step, we use the two types of path automata to test for copying or rearranging. First, we treat copying.

**LEMMA 4.9.** Given a uniform tree transducer  $T$  and an NTA  $N$ , it is decidable in PTIME whether  $T$  is copying over  $\mathcal{L}(N)$ .

**PROOF SKETCH.** We define an NFA  $M$  that accepts text paths satisfying condition (1) or condition (2) of Lemma 4.5. So,  $M$  is non-empty iff  $T$  is copying over  $\mathcal{L}(N)$ . As  $M$  can be constructed in PTIME and testing non-emptiness of NFAs is in PTIME as well, the result follows. Basically,  $M$  simulates both  $A_N$  (as the text path should come from a tree accepted by  $N$ ) and two copies of  $A_T$  (as there should be a path run of  $T$ ) while checking for the existence of two different path runs or while trying to reach a rule in  $T$  which copies.  $\square$

**LEMMA 4.10.** Given a uniform tree transducer  $T$  and an NTA  $N$ , it is decidable in PTIME whether  $T$  is rearranging over  $\mathcal{L}(N)$ .

**PROOF SKETCH.** We will construct in PTIME an NTA  $M$  accepting the set of trees on which  $T$  is rearranging. Therefore, the intersection of  $M$  and  $N$  is non-empty if and only if  $T$  is rearranging over  $\mathcal{L}(N)$ . As the latter can be done in PTIME, the result follows. Intuitively,  $M$  directly searches for nodes satisfying the properties mentioned in Lemma 4.6. We describe  $M$  as if operating in a top-down fashion. Starting from the root of an input tree  $t$ ,  $M$  guesses a path, simulating  $A_T$ . At a certain non-deterministically chosen

point,  $M$  decides it has found state  $q_{i-1}^1 = q_{i-1}^2$  and non-deterministically picks  $q_i^1$  and  $q_i^2$  (cf. Lemma 4.6). From then on,  $M$  continues until it (non-deterministically) arrives at the lowest common ancestor of  $v_1$  and  $v_2$ . On this path,  $M$  will simulate two copies of  $A_T$ ; one copy continuing from state  $q_i^1$  and one copy continuing from state  $q_i^2$ . At the lowest common ancestor of  $v_1$  and  $v_2$ , the simulation of these two copies of  $M$  will split again. Towards  $v_1$ ,  $M$  simulates the copy that started from  $q_i^1$ , and towards  $v_2$ ,  $M$  simulates the copy that started from  $q_i^2$ . Finally,  $M$  accepts when the copy for  $q_i^1$  leads to acceptance of  $A_T$  at node  $v_1$  and the copy of  $q_i^2$  leads to acceptance of  $A_T$  at node  $v_2$ . Recall that the acceptance condition of  $A_T$  means that the text value of the last node is copied to the output of  $T$ .  $\square$

By the previous lemmas the main complexity result of this section readily follows:

**THEOREM 4.11.** *Given a uniform tree transducer  $T$  and an NTA  $N$ , it is decidable in PTIME whether  $T$  is Text-preserving over  $\mathcal{L}(N)$ .*

## 5. DTL

### 5.1 General DTL Transducers

In this section, we consider the abstraction of XSLT called DTL, which was introduced in [11]. DTL programs can navigate in a tree through the use of *binary patterns*. While XSLT employs XPath, we leave the concrete pattern language implicit for now and introduce the following terminology. A *unary pattern* over  $\Sigma$  is a subset of  $\bigcup_{t \in \text{Trees}_\Sigma} (\{t\} \times \text{Nodes}^t)$  and a *binary pattern* over  $\Sigma$  is a subset of  $\bigcup_{t \in \text{Trees}_\Sigma} (\{t\} \times \text{Nodes}^t \times \text{Nodes}^t)$ . We refer to the set of unary and binary patterns over  $\Sigma$  as  $\mathcal{UP}(\Sigma)$  and  $\mathcal{BP}(\Sigma)$ , respectively. We will denote unary patterns by  $\varphi, \psi$  and binary patterns by  $\alpha, \beta$ . To emphasize that unary and binary patterns will be specified by pattern languages, we will also write  $(t, u) \in \varphi$  as  $t \models \varphi(u)$  and  $(t, u, v) \in \alpha$  as  $t \models \alpha(u, v)$ .

Our definition of DTL differs slightly from the one in [11]. To simplify presentation, we disregard construction functions (cf. Section 6 for a discussion). On the other hand, we do extend DTL to deal with text values.

**DEFINITION 5.1.** A DTL *transducer* is a tuple  $T = (\Sigma, \Delta, Q, q_0, R_\Sigma, R_{\text{Text}})$ , where

- $\Sigma$  is the finite alphabet of input symbols,
- $\Delta$  is the finite alphabet of output symbols,
- $Q$  is the finite set of states,
- $q_0 \in Q$  is the initial state,
- $R_\Sigma$  is a finite set of rules of the form  $(q, \varphi) \rightarrow h$ , where  $q \in Q$ ,  $\varphi \in \mathcal{UP}(\Sigma)$ , and  $h$  is in  $\text{Hedges}_\Delta(Q \times \mathcal{BP}(\Sigma))$ . If  $q = q_0$ ,  $h$  is required to be a tree such that  $\text{lab}^h(1) \notin Q \times \mathcal{BP}(\Sigma)$ .<sup>3</sup>
- $R_{\text{Text}}$  is a finite set of rules of the form  $(q, \text{text}) \rightarrow \text{text}$ , with  $q \in Q$ .

To ensure determinism, we require that when  $(q, \varphi) \rightarrow h$  and  $(q, \varphi') \rightarrow h'$  are rules in  $R_\Sigma$ , then there exists no tree  $t$  such

<sup>3</sup>This is just a technical restriction which forces the transducer to output trees.

that  $t \models \varphi(v)$  and  $t \models \varphi'(v)$  for any node  $v$  in  $t$ .<sup>4</sup> Notice in the above definition that there are no restrictions on the allowed patterns. Later, we will use concrete pattern languages like monadic second-order logic and XPath to define patterns.

We are now ready to define the transformation induced by  $T$  on a tree  $t$ . Intuitively,  $T$  starts in state  $q_0$  at the root of  $t$  (note that the root of a tree is always 1). The latter is encoded by the initial configuration  $(q_0, 1)$ . Formally, a *configuration* is a pair in  $Q \times \text{Nodes}^t$ . During computation  $T$  will rewrite the initial configuration into a partial output tree  $\xi \in \text{Trees}_\Delta(Q \times \text{Nodes}^t)$ . That is, a  $\Delta$ -tree where some leaves will be labeled with configurations. The transducer then proceeds by rewriting these configurations extending  $\xi$  until all configurations are gone and the transduction stops. Basically, rewriting of a configuration  $(q, v)$  can be done by a rule  $(q, \varphi) \rightarrow h$  when  $t \models \varphi(v)$ . The output tree is then appended with the hedge  $h$  where the leaves carrying binary patterns result in new configurations.

We next formally define the transduction relation. Thereto, given a tree  $t$ , the transformation relation induced by  $T$  on  $t$ , denoted by  $\Rightarrow_{T,t}$ , is the binary relation on  $\text{Trees}_\Delta(Q \times \text{Nodes}^t)$  defined as follows. For  $\xi, \xi' \in \text{Trees}_\Delta(Q \times \text{Nodes}^t)$ ,  $\xi \Rightarrow_{T,t} \xi'$ , if there is a leaf node  $u \in \text{Nodes}^\xi$  such that

1.  $\text{lab}^\xi(u) = (q, v)$  with  $q \in Q$  and  $v \in \text{Nodes}^t$ ,
2. if  $\text{lab}^t(v) \in \text{Text}$  then  $\xi' = \xi[u \leftarrow \text{lab}^t(v)]$  if there is a rule  $(q, \text{text}) \rightarrow \text{text}$  and  $\xi' = \xi[u \leftarrow \varepsilon]$  otherwise.<sup>5</sup>
3. if  $\text{lab}^t(v) \in \Sigma$ ,
  - if there is a rule  $r = ((q, \varphi) \rightarrow h) \in R_\Sigma$  such that  $t \models \varphi(v)$ , then  $\xi' = \xi[u \leftarrow h\Theta]$ , where  $\Theta$  denotes the substitution of replacing every pair  $(q', \alpha)$  by the hedge  $(q', v_1) \cdots (q', v_m)$ , where
    - $\{v_1, \dots, v_m\} = \{u \mid t \models \alpha(v, u)\}$ , and
    - $v_1 <_{\text{lex}} \cdots <_{\text{lex}} v_m$ ,
  - if there is no rule  $r = ((q, \varphi) \rightarrow h) \in R_\Sigma$  s.t.  $t \models \varphi(v)$ , then  $\xi' = \xi[u \leftarrow \varepsilon]$ .

For any tree  $t$  and DTL transducer  $T$ ,  $T(t)$  is defined to be the tree  $t' \in \text{Trees}_\Delta$ , such that  $(q_0, 1) \Rightarrow_{T,t}^* t'$ , when such a tree exists and is undefined otherwise.<sup>6</sup> As configurations for which no rule applies are rewritten into the empty hedge,  $T(t)$  can only be undefined when  $T$  does not stop. Furthermore, when  $T(t)$  is defined, it is unique because of the determinism restriction on unary patterns.

Example 5.15 contains a DTL transducer which selects all recipes with their description, ingredients, and instructions if they have at least three positive comments. The unary and binary patterns in the example are given in an XPath syntax, formally defined in Section 5.4.

Notice that each top-down tree transducer  $T$  as introduced in the previous section can be defined by a DTL program  $T'$ . Indeed, for each rule  $(q, a) \rightarrow h$  of  $T$ ,  $T'$  has a rule  $(q, \varphi) \rightarrow h'$  where  $t \models \varphi(v) \Leftrightarrow \text{lab}^t(v) = a$ ; and  $h'$  is obtained from  $h$  by replacing each  $q$  that appears as a leaf with the pair  $(q, \text{children})$  where  $\text{children}$  is the pattern selecting all children of the node at hand.

<sup>4</sup>Note that for the tree patterns we consider, testing determinism is decidable.

<sup>5</sup>Recall that  $\varepsilon$  is the empty hedge.

<sup>6</sup>Here,  $R^*$  denotes the reflexive transitive closure of the binary relation  $R$ .

LEMMA 5.2. DTL transducers are admissible.

## 5.2 Copying and Rearranging

Next, we will give an operational characterization of text-preserving DTL transducers in terms of copying and rearranging. Copying and rearranging for DTL transducers is defined analogously as before.

DEFINITION 5.3. For any DTL transducer  $T$  and tree language  $L$ ,  $T$  is *copying over  $L$*  if the transduction defined by  $T$  is copying over  $L$ . It is *rearranging over  $L$* , if the transduction defined by  $T$  is rearranging over  $L$ .

For characterizing these notions, we need some terminology. For the remainder of the discussion, fix a tree  $t$ . We next define when a configuration  $(q', v')$  can follow a configuration  $(q, v)$  in the transduction of  $t$ . Formally, we define  $(q, v) \rightsquigarrow (q', v')$  to hold whenever there are  $\xi, \xi' \in \text{Trees}_\Delta(Q \times \text{Nodes}^t)$  with  $(q_0, 1) \Rightarrow_{T, t}^* \xi \Rightarrow_{T, t} \xi'$  such that  $\text{frontier}(\xi) = \theta_1(q, v)\theta_2$ ,  $\text{frontier}(\xi') = \theta_1\theta_2$ , and  $(q', v') \prec \theta$  for  $\theta_1, \theta_2, \theta \in (\Delta \cup (Q \times \text{Nodes}^t))^*$ . Basically, this says that  $(q, v)$  is rewritten into a tree containing the configuration  $(q', v')$  at its frontier.

A *path run of  $T$  over  $t$*  is a sequence  $(q_0, v_0) \cdots (q_n, v_n)$ , where  $v_0 = 1$  (the root of  $t$ ), and, for all  $i \in [1, n]$ ,  $(q_i, v_i)$  follows  $(q_{i-1}, v_{i-1})$ . A *text path run* is a path run such that  $\text{lab}^t(v_n) \in \text{Text}$  and  $(q_n, \text{text}) \rightarrow \text{text} \in R_{\text{Text}}$ . We say that the path run ends in node  $v_n$ . Intuitively, a text path run is part of a run (which incidentally is a path) of  $T$  on  $t$  which outputs a text value. The next definition determines when a transducer copies, that is, outputs the Text-value carried by the same leaf node  $v$ , either using two different text path runs (ending in  $v$ ), or by using one text path run (ending in  $v$ ) which actually occurs twice in the transduction.

We need one additional notion. Let  $(q, v) \rightsquigarrow (q', v')$ . Then, we say that  $(q', v')$  *doubles* at  $(q, v)$  if there is a rule  $(q, \varphi) \rightarrow h$  and binary patterns  $\alpha_1, \alpha_2 \in \mathcal{BP}(\Sigma)$ , such that  $t \models \varphi(v)$ ,  $t \models \alpha_1(v, v')$ ,  $t \models \alpha_2(v, v')$  and  $(q', \alpha_1) \cdot (q', \alpha_2) \prec \text{frontier}(h)$ . That is, rewriting  $(q, v)$  introduces two occurrences of the configuration  $(q', v')$ .

The following lemma gives our operational characterization of copying for DTL transducers.

LEMMA 5.4. A DTL transducer  $T$  is copying over a tree language  $L$ , if and only if there exists a tree  $t \in L$ , such that there are two different text path runs over  $t$  ending in the same node, or one text path run  $(q_0, v_0) \cdots (q_n, v_n)$  such that  $(q_i, v_i)$  doubles at  $(q_{i-1}, v_{i-1})$  for some  $i \in [1, n]$ .

Next, we characterize when a transducer rearranges two Text-values in a tree. Recall that this happens when there are two leaves  $v$  and  $u$  in  $t$  both carrying a Text-value and where  $v <_{\text{lex}} u$ , but  $T(t)$  swaps the order of the Text-values corresponding to  $v$  and  $u$ . We will show that the latter happens only when during the transduction a rule with rhs  $h$  is used such that

- $(q_u, \alpha_1) \cdot (q_v, \alpha_2) \prec \text{frontier}(h)$ , and the configurations  $(q_u, \alpha_1)$  and  $(q_v, \alpha_2)$  proceed to output the text value at  $u$  and  $v$ , respectively; or,
- when  $(q', \alpha') \prec \text{frontier}(h)$  and  $\alpha'$  selects  $u_1, v_1$  with  $u_1 <_{\text{lex}} v_1$  where  $(q', u_1)$  and  $(q', v_1)$  proceed to output the Text-value at  $u$  and  $v$ , respectively.

We are now ready to prove a characterization for when a DTL transducer is rearranging.

LEMMA 5.5. A DTL transducer  $T$  is rearranging over a tree language  $L$ , if and only if there exists a tree  $t \in L$ , two text path runs  $\theta(q, v)(q_1, v_1) \cdots (q_n, v_n)$  and  $\theta(q, v)(p_1, u_1) \cdots (p_m, u_m)$  with  $v_n <_{\text{lex}} u_m$ , a rule  $(q, \varphi) \rightarrow h$ , and binary patterns  $\alpha_1, \alpha_2$ , such that  $t \models \varphi(v)$ ,  $t \models \alpha_1(v, v_1)$ ,  $t \models \alpha_2(v, u_1)$  and either

- (1)  $(p_1, \alpha_2) \cdot (q_1, \alpha_1) \prec \text{frontier}(h)$ , or
- (2)  $(q_1, \alpha_1) = (p_1, \alpha_2)$ ,  $(q_1, \alpha_1) \prec \text{frontier}(h)$  and  $u_1 <_{\text{lex}} v_1$ .

From Lemmas 5.2, 5.4, and 5.5 and from Theorem 3.3 we now obtain a characterization of text-preserving DTL transducers that can be syntactically tested.

THEOREM 5.6. A DTL transducer  $T$  is not text-preserving over a tree language  $L$  if and only if the condition in Lemma 5.4 or the condition in Lemma 5.5 holds.

## 5.3 Transducers with MSO Transitions

In the present section, we consider the instantiation of DTL with MSO-definable patterns, called  $\text{DTL}^{\text{MSO}}$ . The main result of this section is a proof that testing whether a  $\text{DTL}^{\text{MSO}}$  transduction is text-preserving with respect to a regular tree language, is decidable.

We first define trees over the alphabet  $\Sigma$  as relational structures over the vocabulary  $E, <$  and  $(\text{lab}_\sigma)_{\sigma \in \Sigma}$ . The domain of a tree then is its set of nodes. Furthermore, for two nodes  $v_1, v_2$  of a tree  $t$ ,  $E(v_1, v_2)$  if  $v_2$  is a child of  $v_1$ , and  $v_1 < v_2$  if  $v_1$  and  $v_2$  share the same parent  $v$ , and  $v_1 <_{\text{lex}} v_2$ . Finally, for each  $\sigma \in \Sigma$ ,  $\text{lab}_\sigma$  is the set of nodes labelled with  $\sigma$ . *Monadic Second-order Logic* (MSO) then is the extension of First-order Logic (FO) that allows the use of set variables ranging over subsets of the domain, in addition to the individual variables ranging over elements of the domain, as provided by first-order logic. We assume standard practice. We denote by  $\phi(x_1, \dots, x_k)$  that  $\phi$  is a  $k$ -ary formula with  $k$  free first-order variables. Furthermore, for a tree  $t$ , and nodes  $v_1, \dots, v_k$ , we denote by  $t \models \phi(v_1, \dots, v_k)$  that  $\phi$  holds when the free variables of  $\phi$  are instantiated by  $v_1, \dots, v_k$ . We refer the reader to, e.g., [8] for more background.

Unary and binary patterns will now be defined by unary and binary MSO formulas, respectively. For instance, the formula  $\varphi(x)$  defines the pattern  $\bigcup_{t \in \text{Trees}_\Sigma} \{t\} \times \{u \mid t \models \varphi(u)\}$ . Then,  $\text{DTL}^{\text{MSO}}$  is the instantiation of DTL with MSO-definable patterns. In the definition of  $\text{DTL}^{\text{MSO}}$  transducers, we will simply specify the formulas rather than the patterns they define.

To prove the main result of this section, we need a number of different automata which we define next: tree-jumping automata with MSO transitions, and tree-walking automata with and without MSO tests.

DEFINITION 5.7. A nondeterministic *tree-jumping automaton with MSO transitions* over an alphabet  $\Sigma$ , denoted by  $\text{TJA}^{\text{MSO}}$ , is a tuple  $B = (Q, \Sigma, \delta, q_0, F, M_u, M_b)$  where  $Q$  is a finite set of states,  $q_0 \in Q$  is the initial state,  $F \subseteq Q$  is a set of final states,  $M_u$  is a finite set of unary MSO formulas,  $M_b$  is a finite set of binary MSO formulas, and  $\delta : Q \times M_u \times M_b \rightarrow 2^Q$  is the transition function. A *run* of  $B$  on a tree  $t$  that starts at a node  $v_0$  of  $t$  and ends at a node  $v_n$  of  $t$ , is a sequence  $(q_0, v_0, \varphi_0, \alpha_0) \cdots (q_n, v_n, \varphi_n, \alpha_n)$ , where for each  $i \in [0, n-1]$ ,  $t \models \varphi_i(v_i)$ ,  $t \models \alpha_i(v_i, v_{i+1})$ , and  $q_{i+1} \in \delta(q_i, \varphi_i, \alpha_i)$ . A run is *accepting* if  $q_n \in F$ .

The tree language accepted by  $B$ , denoted  $\mathcal{L}(B)$ , consists of those trees  $t$  for which there exists an accepting run starting



at the root of  $t$ . If there exists a run of  $B$  that starts at a node  $v_1$  in state  $q_1$  and ends at node  $v_2$  in state  $q_2$ , we sometimes also say that, when started in state  $q_1$  in  $v_1$ ,  $B$  ends in state  $q_2$  in  $v_2$ .

We define two restrictions of tree-jumping automata. A nondeterministic *tree-walking automaton with MSO tests*, denoted by  $\text{TWA}^{\text{MSO}}$ , is a  $\text{TJA}^{\text{MSO}}$  where  $M_b$  contains only the predicates  $\text{first-child}(x, y)$ ,  $\text{next-sibling}(x, y)$ ,  $\text{parent}(x, y)$ ,  $\text{previous-sibling}(x, y)$  and equality. Here, the just mentioned predicates have the obvious semantics. For instance,  $\text{first-child}(x, y)$  holds when  $y$  is the first-child of  $x$ , and so on. Note that these are easily MSO-definable.

We show that  $\text{TJA}^{\text{MSO}}$  coincides with the class of unranked regular tree languages strongly bearing on a result by Bloem and Engelfriet [4].

LEMMA 5.8. (1) For each MSO formula  $\alpha(x, y)$ , there exists a  $\text{TWA}^{\text{MSO}}$   $A_\alpha = (Q, \Sigma, \delta, q_0, f, M_u)$ , such that for any tree  $t$  and nodes  $v, u \in \text{Nodes}^t$ ,  $t \models \alpha(v, u)$  if and only if there exists a run of  $A_\alpha$  on  $t$  starting at node  $v$  at state  $q_0$  and ending at node  $u$  at state  $f$ . The converse statement also holds.

(2) For each  $\text{TJA}^{\text{MSO}}$   $B = (Q, \Sigma, \delta, q_0, F, M_u, M_b)$  there exists a  $\text{TWA}^{\text{MSO}}$   $A$ , such that  $\mathcal{L}(B) = \mathcal{L}(A)$ .

PROOF SKETCH. (1) Bloem and Engelfriet [4] have shown the exact same statement for ranked trees. As the first-child next-sibling encoding between ranked and unranked trees is MSO-definable, the result follows for unranked trees as well.

(2) Essentially, the  $\text{TWA}^{\text{MSO}}$   $A$  simulates the  $\text{TJA}^{\text{MSO}}$   $B$  by invoking the automaton  $A_\alpha$  for every transition  $\alpha$  of  $B$ .  $\square$

We note that the constructions in Lemma 5.8 are effective. Furthermore, note that Lemma 5.8(1) together with the fact that MSO-definable unranked tree languages are regular, implies that  $\text{TWA}^{\text{MSO}}$  define the unranked regular tree languages. Therefore, we have the following corollary:

COROLLARY 5.9.  $\text{TJA}^{\text{MSO}}$  define precisely the unranked regular tree languages.

The above implies that emptiness of  $\text{TJA}^{\text{MSO}}$  is decidable. In the remainder of this section we will reduce testing whether a  $\text{DTL}^{\text{MSO}}$  transducer is text-preserving to the emptiness test for  $\text{TJA}^{\text{MSO}}$ . More specifically, for a transducer  $T$ , we will create a  $\text{TJA}^{\text{MSO}}$   $A_T^{\text{copy}}$  and  $A_T^{\text{rearrange}}$  such that  $A_T^{\text{copy}}$  and  $A_T^{\text{rearrange}}$  are empty iff  $T$  is not copying and not rearranging, respectively. By Theorem 3.3,  $T$  then is text-preserving.

We need to introduce a number of automata. We fix a transducer  $T = (Q, \Sigma, \Delta, q_0, R_\Sigma, R_{\text{text}})$ . For states  $q, q' \in Q$ , let  $A_T^{q, q'}$  be the  $\text{TJA}^{\text{MSO}}$  which when started in state  $q$  on node  $v$  of a tree  $t$  ends in state  $q'$  at node  $v'$  iff  $(q, v) \rightsquigarrow^* (q', v')$ . That is, when  $T$  reaches configuration  $(q, v)$  when processing  $t$ ,  $T$  can reach configuration  $(q', v')$ . By  $A_{T, \sigma, \sigma'}^{q, q'}$ , we denote the automaton  $A_T^{q, q'}$  which additionally checks that the start position is labelled with  $\sigma$  and the end position is labelled with  $\sigma'$ .

We start with  $A_T^{\text{rearrange}}$ . The automaton will accept a tree  $t$  when  $T$  is rearranging over  $\{t\}$ . Intuitively, by Lemma 5.5,

the automaton needs to look for the following kind of configurations (ignore the bullets for now):

$$(q_0, 1) \rightsquigarrow^* (q, v) \rightsquigarrow^* \begin{array}{c} \bullet \\ (q_1, v_1) \\ \bullet_1 \end{array} \rightsquigarrow^* \begin{array}{c} (q_n, v_n) \\ \circ_1 \end{array} \rightsquigarrow^* \begin{array}{c} \bullet_2 \\ (p_1, u_1) \\ \circ_2 \end{array} \rightsquigarrow^* \begin{array}{c} (p_m, u_m) \\ \circ_2 \end{array} \quad (\dagger)$$

which satisfy some additional criteria. In particular, we will make sure later that  $v_n$  and  $u_m$  are text nodes and that  $(q, v) \rightsquigarrow (q_1, v_1)$  and  $(q, v) \rightsquigarrow (p_1, u_1)$ , i.e., we can make these transitions in one step.

The automaton makes use of an extended alphabet. Let  $\Sigma_{\text{mark}} = (\Sigma \cup \{\text{text}\}) \times 2^{\{\circ, \circ_1, \circ_2, \bullet, \bullet_1, \bullet_2\}}$ . Intuitively, the alphabet allows us to mark a node with a set of markers. Then, for a marker  $c \in \{\circ, \circ_1, \circ_2, \bullet, \bullet_1, \bullet_2\}$ , we slightly abuse notation and say simply that a node is labelled with  $c$ , when actually there is a  $C \subseteq \{\circ, \circ_1, \circ_2, \bullet, \bullet_1, \bullet_2\}$  with  $c \in C$  and a  $\sigma \in \Sigma$  such that it is labelled with  $(\sigma, C)$ . To check  $(\dagger)$ , we assume  $v$  is labelled with  $\bullet$ ,  $v_1$  is labelled with  $\bullet_1$ ,  $u_1$  is labelled with  $\bullet_2$ ,  $v_n$  is labelled with  $\circ_1$ ,  $u_m$  is labelled with  $\circ_2$ . Then, the configurations in  $(\dagger)$  occur in a tree  $t$  if and only if there is a marking  $t'$  of  $t$  such that

$$t' \in A_{T, \text{root}, \bullet}^{q_0, q} \cap A_{T, \bullet, \bullet_1}^{q, q_1} \cap A_{T, \bullet, \bullet_2}^{q, p_1} \cap A_{T, \bullet_1, \circ_1}^{q_1, q_n} \cap A_{T, \bullet_2, \circ_2}^{p_1, p_m}.$$

Here, the automaton  $A_{T, \text{root}, \bullet}^{q_0, q}$  tests whether the state  $q$  is reachable from  $q_0$  in  $T$ . By  $A_{(\dagger)}^{q, q_1, p_1}$  we refer to this automaton that tests whether there are states  $q_n, p_m$  such that condition  $(\dagger)$  holds for given states  $q, q_1, p_1$ .

We need some additional regular languages over  $\Sigma_{\text{mark}}$ . They can for instance easily be defined in MSO. We assume that they are represented by NTAs:

- (i)  $A_{\text{mark}}$ : for all  $c \in \{\bullet, \bullet_1, \bullet_2, \circ_1, \circ_2\}$ , there is exactly one node labelled with  $c$ . Furthermore, the node labelled with  $\circ_1$  (resp.,  $\circ_2$ ), is a leaf labelled with  $(\text{text}, C)$  and  $\circ_1 \in C$  (resp.,  $\circ_2 \in C$ ). We will abuse notation and refer to the nodes labelled with  $c$  simply as  $c$ .
- (ii)  $A_{<, \bullet}$ :  $\bullet_2 <_{\text{lex}} \bullet_1$ ;
- (iii)  $A_{<, \circ}$ :  $\circ_1 <_{\text{lex}} \circ_2$ ;
- (iv)  $A_{\bullet}^\varphi$ :  $t \models \varphi(\bullet)$ ;
- (v)  $A_{\bullet_1}^\alpha$ :  $t \models \alpha(\bullet, \bullet_1)$ ; and,
- (vi)  $A_{\bullet_2}^\beta$ :  $t \models \beta(\bullet, \bullet_2)$ .

We define two automata  $A_1^{\text{rearrange}}$  and  $A_2^{\text{rearrange}}$  corresponding to the two conditions of Lemma 5.5. To state the conditions, we assume as before that the nodes are marked as in condition  $(\dagger)$ .

In the first condition, we need to test, in addition to  $(\dagger)$ , whether  $\circ_1 <_{\text{lex}} \circ_2$  and there is a rule  $(q, \varphi) \rightarrow h$  in  $T$  with  $(p_1, \beta) \cdot (q_1, \alpha) \prec \text{frontier}(h)$  such that  $t \models \varphi(\bullet)$ ,  $t \models \alpha(\bullet, \bullet_1)$ , and  $t \models \beta(\bullet, \bullet_2)$ . In order to formally capture this condition, we define a relation  $G$  that contains precisely the tuples  $(q, \varphi, p_1, \beta, q_1, \alpha)$  such that there exists a rule  $(q, \varphi) \rightarrow h$  in  $T$  with  $(p_1, \beta) \cdot (q_1, \alpha) \prec \text{frontier}(h)$ . We are now ready to define  $A_1^{\text{rearrange}}$ :

$$A_1^{\text{rearrange}} := A_{\text{mark}} \cap A_{<, \circ} \cap \left( \bigcup_{(q, \varphi, p_1, \beta, q_1, \alpha) \in G} (A_{(\dagger)}^{q, q_1, p_1} \cap A_{\bullet}^\varphi \cap A_{\bullet_1}^\alpha \cap A_{\bullet_2}^\beta) \right)$$

In the second condition of Lemma 5.5, we need to test, in addition to  $(\dagger)$ , whether  $\circ_1 <_{\text{lex}} \circ_2$  and there is a rule  $(q, \varphi) \rightarrow h$  in  $T$  with  $(q_1, \alpha) \prec \text{frontier}(h)$  such that  $(q_1, \alpha) = (p_1, \beta)$ ,  $\bullet_2 <_{\text{lex}} \bullet_1$ ,  $t \models \varphi(\bullet)$ ,  $t \models \alpha(\bullet, \bullet_1)$ , and  $t \models \alpha(\bullet, \bullet_2)$ .

In order to formally capture this condition, we define a relation  $H$  that contains precisely the tuples  $(q, \varphi, q_1, \alpha)$  such that there exists a rule  $(q, \varphi) \rightarrow h$  in  $T$  with  $(q_1, \alpha) \prec \text{frontier}(h)$ . We are now ready to define  $A_2^{\text{rearrange}}$ :

$$A_2^{\text{rearrange}} := A_{\text{mark}} \cap A_{<, \circ} \cap A_{<, \bullet} \left( \bigcup_{(q, \varphi, q_1, \alpha) \in H} (A_{(\ddagger)}^{q, q_1, q_1} \cap A_{\bullet}^{\varphi} \cap A_{\bullet_1}^{\alpha} \cap A_{\bullet_2}^{\alpha}) \right)$$

LEMMA 5.10. *Given a DTL<sup>MSO</sup> transducer  $T$ , the language defined by the TJA<sup>MSO</sup>  $A_1^{\text{rearrange}} \cup A_2^{\text{rearrange}}$  is non-empty if and only if  $T$  is rearranging over  $\text{Trees}_{\Sigma}$ .*

The overall procedure for  $A_T^{\text{copy}}$  is similar. The automaton will accept a tree  $t$  when  $T$  is copying over  $\{t\}$ . Intuitively, by Lemma 5.4, the automaton needs to look for the following kind of configurations

$$(q_0, 1) \rightsquigarrow^* (q, u) \rightsquigarrow^* \begin{array}{ccc} (q_1, u_1) & \rightsquigarrow^* & (q_n, v) \\ \bullet_1 & & \circ \\ \bullet_2 & & \circ \\ (p_1, u_2) & \rightsquigarrow^* & (p_m, v) \end{array} \quad (\ddagger)$$

which satisfy some additional criteria. The node  $v$  should be a text node and the transitions from  $(q, u)$  to  $(q_1, u_1)$  and to  $(p_1, u_2)$  should be possible in one step. Note that both paths in the picture lead to the same node  $v$ .

As before the automaton makes use of the extended alphabet  $\Sigma_{\text{mark}}$ . With the same conventions as before, to check  $(\ddagger)$ , we assume  $u$  is labelled with  $\bullet$ ,  $u_1$  is labelled with  $\bullet_1$ ,  $u_2$  is labelled with  $\bullet_2$ , and  $v$  is labelled with  $\circ$ . Then, the configurations in  $(\ddagger)$  occur in a tree  $t$  if and only if there is a marking  $t'$  of  $t$  such that

$$t' \in A_{T, \text{root}, \bullet}^{q_0, q} \cap A_{T, \bullet, \bullet_1}^{q, q_1} \cap A_{T, \bullet, \bullet_2}^{q, p_1} \cap A_{T, \bullet_1, \circ}^{q_1, q_n} \cap A_{T, \bullet_2, \circ}^{p_1, p_m}$$

Analogously as before, the automaton  $A_{T, \text{root}, \bullet}^{q_0, q}$  tests whether the state  $q$  is reachable from  $q_0$  in  $T$ . By  $A_{(\ddagger)}^{q, q_1, p_1}$  we refer to this automaton that tests whether there are states  $q_n, p_m$  such that condition  $(\ddagger)$  holds for given states  $q, q_1, p_1$ .

In addition to the NTAs for the regular languages (i)–(vi) above, we need

(vii)  $A'_{\text{mark}}$ : for any  $c \in \{\bullet, \bullet_1, \bullet_2, \circ\}$ , there is exactly one node labelled with  $c$  and the node labelled  $\circ$  is a text node in  $t$ ,

(viii)  $A_{\bullet_1 \neq \bullet_2}$ :  $\bullet_1 \neq \bullet_2$ .

(ix)  $A_{\bullet_1 = \bullet_2}$ :  $\bullet_1 = \bullet_2$ .

We define two automata  $A_1^{\text{copy}}$  and  $A_2^{\text{copy}}$  corresponding to the two conditions of Lemma 5.4. To state the conditions, we assume as before that the nodes are marked as in condition  $(\ddagger)$ . We also assume that the nodes labelled with  $\bullet, \bullet_1, \bullet_2, \circ$  are unique.

In the first condition, we need to test, in addition to  $(\ddagger)$ , if there are two different text path runs over  $t$  ending in the same node  $v$ . Since the path from  $(q_0, 1)$  to  $(q, u)$  can have arbitrary length, we can assume w.l.o.g. that the last configuration in which these two text path runs are the same is  $(q, u)$ . This means that there exists a rule  $(q, \varphi) \rightarrow h$  in  $T$  with  $(q_1, \alpha) \prec \text{frontier}(h)$  and  $(p_1, \beta) \prec \text{frontier}(h)$  such that  $t \models \varphi(\bullet)$ ,  $t \models \alpha(\bullet, \bullet_1)$ , and  $t \models \beta(\bullet, \bullet_2)$ . Furthermore, since  $(q, u)$  is the last configuration for which the text path runs are the same, we either have that (a)  $q_1 \neq p_1$  or (b)  $\bullet_1 \neq \bullet_2$ . In order to formalize this condition, let  $I$  be the set that contains all tuples  $(q, \varphi, q_1, \alpha, p_1, \beta)$  such that  $(q, \varphi) \rightarrow h$  is a

rule in  $T$  with  $(q_1, \alpha) \prec \text{frontier}(h)$  and  $(p_1, \beta) \prec \text{frontier}(h)$ . Then, we can define  $A_{1a}^{\text{copy}}$  and  $A_{1b}^{\text{copy}}$ :

$$A_{1a}^{\text{copy}} := A'_{\text{mark}} \cap \left( \bigcup_{(q, \varphi, p_1, \beta, q_1, \alpha) \in I \wedge p_1 \neq q_1} (A_{(\ddagger)}^{q, q_1, p_1} \cap A_{\bullet}^{\varphi} \cap A_{\bullet_1}^{\alpha} \cap A_{\bullet_2}^{\beta}) \right)$$

$$A_{1b}^{\text{copy}} := A'_{\text{mark}} \cap A_{\bullet_1 \neq \bullet_2} \left( \bigcup_{(q, \varphi, p_1, \beta, q_1, \alpha) \in I} (A_{(\ddagger)}^{q, q_1, p_1} \cap A_{\bullet}^{\varphi} \cap A_{\bullet_1}^{\alpha} \cap A_{\bullet_2}^{\beta}) \right)$$

In the second condition of Lemma 5.4 we need to test whether there exists a text path run with two consecutive configurations such that the second one doubles at the first one. W.l.o.g., we can assume that  $(q_1, u_1)$  doubles at  $(q, u)$ . This means that there exists a rule  $(q, \varphi) \rightarrow h$  in  $T$  with  $(q_1, \alpha) \cdot (q_1, \beta) \prec \text{frontier}(h)$ , such that  $t \models \varphi(\bullet)$ ,  $t \models \alpha(\bullet, \bullet_1)$ ,  $t \models \beta(\bullet, \bullet_2)$ , and  $u_1 = u_2$ , i.e., the node  $u_1$  is labelled by  $\bullet_1$  and  $\bullet_2$ . In order to formalize this condition, let  $J$  be the set that contains all tuples  $(q, \varphi, q_1, \alpha, \beta)$  such that  $(q, \varphi) \rightarrow h$  is a rule in  $T$  with  $(q_1, \alpha) \cdot (q_1, \beta) \prec \text{frontier}(h)$ . Then, we can define  $A_2^{\text{copy}}$ :

$$A_2^{\text{copy}} := A'_{\text{mark}} \cap A_{\bullet_1 = \bullet_2} \cap \left( \bigcup_{(q, \varphi, q_1, \alpha, \beta) \in J} (A_{(\ddagger)}^{q, q_1, q_1} \cap A_{\bullet}^{\varphi} \cap A_{\bullet_1}^{\alpha} \cap A_{\bullet_2}^{\beta}) \right)$$

LEMMA 5.11. *Given a DTL<sup>MSO</sup> transducer  $T$ , the language defined by the TJA<sup>MSO</sup>  $A_{1a}^{\text{copy}} \cup A_{1b}^{\text{copy}} \cup A_2^{\text{copy}}$  is non-empty if and only if  $T$  is copying over  $\text{Trees}_{\Sigma}$ .*

To test whether a DTL<sup>MSO</sup> transducer  $T$  is text-preserving over a regular tree language  $\mathcal{L}$ , it follows from Theorem 3.3 that it is sufficient to test the conditions in Lemmas 5.10 and 5.11. We therefore obtain the following theorem.

THEOREM 5.12. *Given a DTL<sup>MSO</sup> transducer  $T$  and a regular language  $\mathcal{L}$ , it is decidable whether  $T$  is text-preserving.*

As satisfiability of MSO is non-elementary over trees, testing whether a DTL<sup>MSO</sup> transducer is text-preserving is non-elementary as well. Indeed, consider a transducer which outputs two copies of the input tree when an MSO-sentence is satisfied at the root. Then this transduction is text-preserving if and only if the formula is not satisfiable.

## 5.4 Transducers with XPath Transitions

We recall the definition of Core XPath. Since we do not use other fragments of XPath, we will use the term XPath to refer to Core XPath. We use the definition from [24]. We denote XPath node expressions by  $\varphi, \psi, \dots$  and path expressions by  $\alpha, \beta, \dots$

DEFINITION 5.13. XPath node expressions and path expressions are defined by simultaneous induction, as follows: Path expressions:  $\alpha ::= R \mid R^* \mid \cdot \mid \alpha/\beta \mid \alpha \cup \beta \mid \alpha[\varphi]$

Node expressions:  $\varphi ::= \sigma \mid \langle \alpha \rangle \mid \top \mid \neg \varphi \mid \varphi \wedge \psi$

Here,  $\sigma \in \Sigma$  is a label test and  $R$  is one of the relations child, parent, next-sibling, or previous-sibling. We denote these relations by  $\downarrow, \uparrow, \rightarrow,$  and  $\leftarrow$ , respectively.

$\llbracket R \rrbracket_{\text{PEExpr}}$	=	the pairs in relation $R$
$\llbracket R^* \rrbracket_{\text{PEExpr}}$	=	reflexive and transitive closure of $\llbracket R \rrbracket_{\text{PEExpr}}$
$\llbracket \cdot \rrbracket_{\text{PEExpr}}$	=	$\{(u, u) \mid u \in \text{Nodes}^t\}$
$\llbracket \alpha/\beta \rrbracket_{\text{PEExpr}}$	=	$\{(u, w) \mid \exists v. (u, v) \in \llbracket \alpha \rrbracket_{\text{PEExpr}} \text{ and } (v, w) \in \llbracket \beta \rrbracket_{\text{PEExpr}}\}$
$\llbracket \alpha \cup \beta \rrbracket_{\text{PEExpr}}$	=	$\llbracket \alpha \rrbracket_{\text{PEExpr}} \cup \llbracket \beta \rrbracket_{\text{PEExpr}}$
$\llbracket \alpha/\varphi \rrbracket_{\text{PEExpr}}$	=	$\{(u, v) \in \llbracket \alpha \rrbracket_{\text{PEExpr}} \mid v \in \llbracket \varphi \rrbracket_{\text{NExpr}}\}$
$\llbracket \sigma \rrbracket_{\text{NExpr}}$	=	$\{u \in \text{Nodes}^t \mid \text{lab}^t(u) = \sigma\}$
$\llbracket \langle \alpha \rangle \rrbracket_{\text{NExpr}}$	=	$\{u \in \text{Nodes}^t \mid \exists v \in \text{Nodes}^t. (u, v) \in \llbracket \alpha \rrbracket_{\text{PEExpr}}\}$
$\llbracket \top \rrbracket_{\text{NExpr}}$	=	$\text{Nodes}^t$
$\llbracket \neg \varphi \rrbracket_{\text{NExpr}}$	=	$\text{Nodes}^t \setminus \llbracket \varphi \rrbracket_{\text{NExpr}}$
$\llbracket \varphi \wedge \psi \rrbracket_{\text{NExpr}}$	=	$\llbracket \varphi \rrbracket_{\text{NExpr}} \cap \llbracket \psi \rrbracket_{\text{NExpr}}$

**Table 1: Semantics of XPath.**

The semantics of XPath expressions relative to a tree  $t$  is given by the functions  $\llbracket \cdot \rrbracket_{\text{PEExpr}}^t$  and  $\llbracket \cdot \rrbracket_{\text{NExpr}}^t$ . These functions are defined in Table 1. We omit the tree  $t$  if it is clear from the context. In the remainder of this section, we denote by  $\varphi(u)$  and  $\alpha(u, v)$  that  $u \in \llbracket \varphi \rrbracket_{\text{NExpr}}$  and  $(u, v) \in \llbracket \alpha \rrbracket_{\text{PEExpr}}$ , respectively.

**DEFINITION 5.14.** A  $\text{DTL}^{\text{XPath}}$  transducer is a DTL transducer in which all unary and binary patterns are XPath node and XPath path expressions, respectively.

**EXAMPLE 5.15.** The following DTL transducer selects the descriptions, ingredients, and instructions from all recipes that have at least three positive comments in our running example.

$(q_0, \text{recipes})$	$\rightarrow$	$\text{recipes}((q, \downarrow))$
$(q, \varphi)$	$\rightarrow$	$\text{recipe}((q, \downarrow))$
$(q, \sigma)$	$\rightarrow$	$\sigma((q, \downarrow))$ ( $\sigma \in \{\text{description, ingredients, br, instructions}\}$ )
$(q, \text{item})$	$\rightarrow$	$(q, \downarrow)$
$(q, \text{text})$	$\rightarrow$	$\text{text}$

where

$$\varphi = \text{recipe} \wedge \langle \downarrow [\text{comments}] / \downarrow [\text{positive}] / \downarrow [\text{comment}] / \rightarrow [\text{comment}] / \rightarrow [\text{comment}] \rangle.$$

As in Section 5.3, we will reduce testing whether a  $\text{DTL}^{\text{XPath}}$  transducer is text-preserving to the emptiness test of an automaton. The beginning of our construction is similar to the one in Section 5.3. We will simulate partial computations of a  $\text{DTL}^{\text{XPath}}$  transducer by a *tree-jumping automaton with XPath transitions* ( $\text{TJA}^{\text{XPath}}$ ). A  $\text{TJA}^{\text{XPath}}$  is a  $\text{TJA}^{\text{MSO}}$  in which all unary and binary formulas are XPath node and XPath path expressions, respectively.

Fix a  $\text{DTL}^{\text{XPath}}$  transducer  $T = (Q, \Sigma, \Delta, q_0, R_\Sigma, R_{\text{Text}})$ . For states  $q, q' \in Q$  and  $\sigma, \sigma' \in \Sigma$ , we denote by  $A_{T, \sigma, \sigma'}^{q, q'}$  the  $\text{TJA}^{\text{XPath}}$  which, when started in state  $q$  on some node  $v$  in a tree  $t$  ends in state  $q'$  at node  $v'$  if and only if  $(q, v) \rightsquigarrow^* (q', v')$  in  $T$ . Furthermore, it checks whether  $v$  is labelled by  $\sigma$  and  $v'$  is labelled by  $\sigma'$ . Furthermore, each of these automata can be constructed from  $T$  in linear time.

It remains to show that we can perform a similar construction as in Section 5.3. However, in this section, we also have to mind the complexity bounds. In order to avoid a large blow-up when performing unions and intersections of

automata, we will translate  $\text{TJA}^{\text{XPath}}$  to two-way alternating tree walking automata using techniques from [3].

**LEMMA 5.16.** *For each  $\text{TJA}^{\text{XPath}}$ , we can construct an equivalent 2ATWA in polynomial time.*

Two-way alternating tree walking automata have the interesting property that their intersections and unions can be constructed very efficiently. More precisely, for 2ATWAs  $A_1, \dots, A_n$ , we can construct a 2ATWA for  $\mathcal{L}(A_1) \cup \dots \cup \mathcal{L}(A_n)$  and a 2ATWA for  $\mathcal{L}(A_1) \cap \dots \cap \mathcal{L}(A_n)$  in time linear in  $|A_1| + \dots + |A_n|$ . This implies that we can perform the constructions from Section 5.3 on the 2ATWAs in polynomial time.

**LEMMA 5.17.** *Given a  $\text{DTL}^{\text{XPath}}$  transducer  $T$  and an NTA  $N$ , we can construct, in polynomial time, a 2ATWA  $A$  such that  $T$  is text-preserving over  $\mathcal{L}(N)$  if and only if  $\mathcal{L}(A) = \emptyset$ .*

Since testing emptiness of a 2ATWA is in  $\text{EXPTIME}$  [5], we now know that testing whether a  $\text{DTL}^{\text{XPath}}$  transducer is text-preserving is in  $\text{EXPTIME}$  as well. Furthermore, since XPath satisfiability w.r.t. a DTD is  $\text{EXPTIME}$ -complete [17, 21], testing whether a  $\text{DTL}^{\text{XPath}}$  transducer is text-preserving is also  $\text{EXPTIME}$ -hard.

**THEOREM 5.18.** *Given a  $\text{DTL}^{\text{XPath}}$  transducer  $T$  and an NTA  $N$ , deciding whether  $T$  is text-preserving over  $\mathcal{L}(N)$  is  $\text{EXPTIME}$ -complete.*

Actually,  $\text{EXPTIME}$ -hardness can also be obtained through a reduction from emptiness of deterministic tree-walking automata known to be  $\text{EXPTIME}$ -complete [20, 23].

## 6. RELATED WORK

The advent of XML induced a renewed interest in tree transducers. Milo, Suci, and Vianu used  $k$ -pebble transducers in their seminal paper [18] as a general model for XML transformations, encompassing a variety of languages like XML-QL, XQuery, and XSLT. Engelfriet and Maneth showed that these form a subclass of the macro tree transducers [7].

In the present paper, we focus on XSLT-like transformations. The simplest of these are the top-down uniform transducers introduced in [13]. A more accurate formalization capturing the navigational power of XSLT as well, is given by DTL as introduced by Maneth and Neven [11]. We extend DTL to deal with text values, albeit in a modest way in accordance with our focus on simple transformations: when encountering Text-values, they can only be immediately output or discarded. Transducer models for XSLT equipped to employ text- or data values during computation have been considered in [2, 19]. To simplify presentation, we disregarded the construction functions employed in [11]. Intuitively, the difference is as follows. DTL as defined here associates one state with every binary pattern, while through the use of construction functions, the selection mechanism is generalized, by associating several states to every binary pattern. As this level of expressiveness is not present in XSLT, we decided to neglect it.

The most deeply studied problem concerning XML and transducers is undoubtedly the XML typechecking problem [18, 13, 14, 15, 9, 12, 22, 10], where it is asked whether  $T(t) \in S_{\text{out}}$  for all  $t \in S_{\text{in}}$  for given input and output

schema  $S_{in}$  and  $S_{out}$  and transducer  $T$ . The typechecking problem is quite different from the problem considered in this paper, both in statement as in computational difficulty. For instance, typechecking top-down uniform tree transducers against unranked tree automata is already EXPTIME-complete while testing whether one is text-preserving is in PTIME for the corresponding setting.

Finally, we note that our Text-trees are different from the data trees in, e.g., [6], where each node carries a label from  $\Sigma$  and a label from Text. Moreover, unlike the logics and automata described in [6], the presented transducer model can not perform equality tests on Text-values.

## 7. CONCLUSION

Text-preserving transformations can be seen as a well-defined simple class of transformations for text-centric XML documents. In this paper, we showed that for XSLT-style transformations this property is decidable and even tractable for the top-down case.

From the employed proof technique it follows that for every transducer in one of the above mentioned classes, the maximal subset of the input schema can be computed on which the transformation is text-preserving. Indeed, we have actually shown that the class of trees on which a transducer is not text-preserving is a regular tree language. It would be interesting to pinpoint the exact size of a representation of that maximal subclass.

Another interesting question for future research is whether there is a normal form for text-preserving transformations in the considered formalisms.

Finally, although we think that mappings which restrict the output to be a subsequence of the input are relevant, there might be other notions that could be interesting in the context of text-centric mappings. It would therefore be interesting to develop more general notions of which text-preserving is a special case.

Actually, our employed proof technique already allows to test stronger properties than just text-preserving. For instance, we can specify a set of tree types (for example, a set of node labels) under which text values in the subtree can not be deleted. In our running example, we could require, e.g., that the transformation is text-preserving and that it does not delete any text values under a node labelled **instructions**. These more flexible tests do not influence the complexity of the problem.

## Acknowledgment

We thank Maarten Marx for bringing the notion of text-preserving mappings and its relevance for text-centric documents to our attention.

## 8. REFERENCES

- [1] J. Albert, D. Giammerresi, D. Wood. Normal form algorithms for extended context free grammars. *Theor. Comp. Sc.*, 267(1–2):35–47, 2001.
- [2] G. J. Bex, S. Maneth, F. Neven. A formal model for an expressive fragment of XSLT. *Inf. Syst.*, 27(1):21–39, 2002.
- [3] H. Björklund, W. Gelade, W. Martens. Incremental XPath evaluation. *ACM Trans. Database Syst.*, 35(4), 2011.
- [4] R. Bloem, J. Engelfriet. A comparison of tree transductions defined by monadic second order logic and by attribute grammars. *J. Comput. Syst. Sci.*, 61(1):1–50, 2000.
- [5] M. Bojanczyk. Tree-walking automata. In *LATA*, pages 1–2, 2008.
- [6] M. Bojanczyk, A. Muscholl, T. Schwentick, L. Segoufin. Two-variable logic on data trees and XML reasoning. *Journal of the ACM*, 56(3), 2009.
- [7] J. Engelfriet, S. Maneth. A comparison of pebble tree transducers with macro tree transducers. *Acta Inf.*, 39(9):613–698, 2003.
- [8] L. Libkin. *Elements Of Finite Model Theory*. Springer Verlag, 2004.
- [9] S. Maneth, A. Berlea, T. Perst, H. Seidl. XML type checking with macro tree transducers. In *PODS*, pages 283–294, 2005.
- [10] S. Maneth, S. Friese, H. Seidl. Type checking of tree walking transducers. In *Modern Applications of Automata Theory*. World Scientific Publishing, 2011.
- [11] S. Maneth, F. Neven. Structured document transformations based on XSL. In *DBPL*, pages 80–98, 1999.
- [12] S. Maneth, T. Perst, H. Seidl. Exact XML type checking in polynomial time. In *ICDT*, pages 254–268, 2007.
- [13] W. Martens, F. Neven. On the complexity of typechecking top-down XML transformations. *Theor. Comp. Sc.*, 336(1):153–180, 2005.
- [14] W. Martens, F. Neven. Frontiers of tractability for typechecking simple XML transformations. *J. Comput. Syst. Sci.*, 73(3):362–390, 2007.
- [15] W. Martens, F. Neven, M. Gyssens. Typechecking top-down XML transformations: Fixed input or output schemas. *Inf. and Comput.*, 206(7):806–827, 2008.
- [16] W. Martens, F. Neven, T. Schwentick. Complexity of decision problems for XML schemas and chain regular expressions. *SIAM J. Comput.*, 39(4):1486–1530, 2009.
- [17] M. Marx. XPath with conditional axis relations. In *EDBT*, pages 477–494, 2004.
- [18] T. Milo, D. Suci, V. Vianu. Typechecking for XML transformers. *J. Comput. Syst. Sci.*, 66(1):66–97, 2003.
- [19] F. Neven. On the power of walking for querying tree-structured data. In *PODS*, pages 77–84, 2002.
- [20] F. Neven. Attribute grammars for unranked trees as a query language for structured documents. *J. Comput. Syst. Sci.*, 70(2):221–257, 2005.
- [21] F. Neven, T. Schwentick. On the complexity of XPath containment in the presence of disjunction, DTDs, and variables. *Log. Meth. in Comp. Sc.*, 2(3), 2006.
- [22] T. Perst and H. Seidl. Macro forest transducers. *Inf. Process. Lett.*, 89(3):141–149, 2004.
- [23] M. Samuelides and L. Segoufin. Complexity of pebble tree-walking automata. In *FCT*, pages 458–469, 2007.
- [24] B. ten Cate and C. Lutz. The complexity of query containment in expressive fragments of XPath 2.0. *Journal of the ACM*, 56(6), 2009.