

# Equivalence of Symbolic Tree Transducers

Vincent Hugot, Adrien Boiret, Joachim Niehren

► **To cite this version:**

Vincent Hugot, Adrien Boiret, Joachim Niehren. Equivalence of Symbolic Tree Transducers. DLT 2017 - Developments in Language Theory, Aug 2017, Liege, Belgium. pp.12, 10.1007/978-3-642-29709-0\_32. hal-01517919v2

**HAL Id: hal-01517919**

**<https://hal.inria.fr/hal-01517919v2>**

Submitted on 30 Sep 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Equivalence of Symbolic Tree Transducers

Vincent Hugot<sup>2,3</sup>, Adrien Boiret<sup>2,3</sup>, and Joachim Niehren<sup>1,3</sup>

<sup>1</sup> Inria, Lille, France

<sup>2</sup> University of Lille, France

<sup>3</sup> Links project (Inria & Cristal Lab - UMR CNRS 9189)

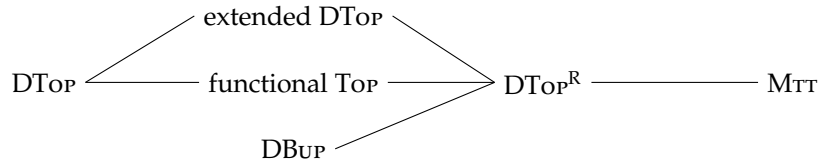
**Abstract.** Symbolic tree transducers are programs that transform data trees with an infinite signature. In this paper, we show that the equivalence problem of deterministic symbolic top-down tree transducers (DTOP) can be reduced to that of classical DTOP. As a consequence the equivalence of two symbolic DTOP can be decided in NEXPTIME, when assuming that all operations related to the processing of data values are in PTIME. This result can be extended to symbolic DTOP with lookahead and thus to deterministic symbolic bottom-up tree transducers.

This is an extended version of the paper published at DLT 2017.

## 1 Introduction

Data trees are widely used in various domains of computer science. They represent programs in compiler construction or program analysis, syntactic sentence structure in computational linguistics, all or part of the database instances in semi-structured databases, and structured documents in document processing. The most widely used current formats for data trees are JSON (the Java Script Object Notation) and XML (the eXtensible Markup Language).

We are interested in deciding the equivalence of programs that define transformations on data trees. For instance, we may consider XSLT programs defining XML transformations or Linux installation scripts written in bash that change the file system tree. Our approach is to compile a subclass of such programs into classes of tree transducers for which equivalence is decidable. Here we present a partial landscape of classical classes of tree transducers without data values [7, 9, 13], where inclusion is read from left to right:



The class DTOP<sup>R</sup> of deterministic top-down tree transducers with regular lookahead by a deterministic bottom-up tree automaton is particularly well behaved [9]. It is closed under composition, which makes it suitable for compilation of programs, and its

equivalence problem is decidable in NEXPTIME, by PTIME reduction to the equivalence problem of the class DTOP, for which equivalence is decidable in NEXPTIME [15, 12, 18]. Furthermore, the class DTOP<sup>R</sup> subsumes three other classes of tree transducers with pairwise incomparable expressiveness, which capture different aspects of programs: extended DTOP with nested pattern matching, functional top-down tree transducers (functional TOP) that relax the determinism requirement of DTOP, and deterministic bottom-up tree transducers DBUP that operate the other way around. The more general class of macro tree transducers (MTT) is much more expressive (and includes lookaheads), but has a long-standing open equivalence problem [10, 11, 19] and fails to be closed under composition (though its linear size increase subclass has better properties).

In contrast to classical machines that operate on ranked trees over finite signatures, what we need for program verification are generalised machines that operate on data trees with infinite signatures. Most typically the data values which label the nodes of data trees may be strings over some finite alphabet or natural numbers. For dealing with data trees, the classical classes of tree transducers were extended to symbolic classes [21, 20, 16], and similarly for other kinds of finite state machines. The general idea is to use patterns for describing infinitely many data values in a finite manner, and to allow the transducers to apply transformations on the data values themselves.

We first illustrate by an example that the class of symbolic extended DTOP is relevant in practice. For this, we consider the following extended DTOP, which performs a routine cleanup and statistics task on a list of log files in a file system, as illustrated in the example of Figure 1.

$$q\langle \text{NIL} \rangle \rightarrow \text{CONS}\langle \text{FILE}\langle \text{"log"}, \text{""} \rangle, \text{NIL} \rangle \quad (1)$$

$$q\langle \text{CONS}\langle x_1, x_2 \rangle \rangle \rightarrow \text{CONS}\langle q_{\text{name}}\langle x_1 \rangle, q\langle x_2 \rangle \rangle \quad (2)$$

$$q_{\text{name}}\langle \text{FILE}\langle \text{"log"}, x_2 \rangle \rangle \rightarrow \text{FILE}\langle \text{"stats.1"}, f_{\text{stats}}@x_2 \rangle \quad (3)$$

$$q_{\text{name}}\langle \text{FILE}\langle x_1 : \text{"stats."}\langle \text{"0"}..\langle \text{"9"} \rangle^+ \rangle, x_2 \rangle \rightarrow \text{FILE}\langle f_{\text{incr}}@x_1, f_{\text{id}}@x_2 \rangle \quad (4)$$

Such transducers have nested patterns with variables  $x_1, x_2, \dots$  for matching subtrees, expressions for matching data values such as  $\text{"stats."}\langle \text{"0"}..\langle \text{"9"} \rangle^+$  or  $\text{CONS}$ , and applications of externally defined functions, such as  $f_{\text{stats}}@x_2$ , where  $f_{\text{stats}}$  is a string

transformation that produces statistics from its input string (log contents). It should be noted that symbolic functional  $\text{Top}$  are insufficient for this example since rules such as (3) and (4), with nested patterns, cannot be expressed in a top-down manner. In contrast, symbolic  $\text{DBUF}$  offer an alternative solution for this concrete example.

Veanes and Bjørner [21, 20] started the study of symbolic transducers. They showed that equivalence is decidable for symbolic functional  $\text{Top}$ , if the corresponding problems on data pattern and transformations are. In this paper, we notice that the landscape of tree transducers above remains unchanged when turning classes of classical tree transducers symbolic. Therefore, we can show that the equivalence problem is decidable for the symbolic counterparts of all classes in the landscape except for  $\text{MTT}$ . To see this, note that any symbolic  $\text{DTOP}$  is a symbolic functional  $\text{Top}$ , so equivalence for symbolic  $\text{DTOP}$  is decidable. Furthermore, equivalence for symbolic  $\text{DTOP}^R$ s can be reduced to equivalence of symbolic  $\text{DTOP}$  as in the classical case.

We then start studying the complexity of equivalence for classes of symbolic tree transducers. Our main result is that equivalence for symbolic  $\text{DTOP}^R$ s is in  $\text{NEXP TIME}$ , under the assumption that operations on patterns and data transformations can be performed in  $\text{P TIME}$ . If not, one needs to multiply the worst case exponential time with the maximal time needed for such operations. We obtain this result from a novel reduction from the equivalence of symbolic  $\text{DTOP}$  to the equivalence of classical  $\text{DTOP}$ , using a weakened version of the *origin-equivalence* in [3]. This reduction allows us to conclude that the equivalence problem of symbolic  $\text{DTOP}$  is indeed in  $\text{NEXP TIME}$  as for classical  $\text{DTOP}$  (and not in  $2\text{NEXP TIME}$  as a naive analysis would lead to believe). Due to the modularity of the construction, the equivalence testers obtained for  $\text{DTOP}^R$ s are easy to prove correct, to analyse, and to implement.

## 2 Tree Automata and Transducers

Some familiarity with formal languages and automata theory, as covered for instance in [5], is assumed.

Given a set  $S$ , we denote its cardinality by  $|S|$  and its powerset by  $2^S$ . The set of Boolean values is written  $\mathbb{B} = \{0, 1\}$ .  $\mathbb{N}$  is the set of natural integers, including zero. We write  $m..n$  the integer interval  $[m, n] \cap \mathbb{N}$ . We will denote tuples  $(a_0, a_1, \dots, a_n)$  by  $a_0 \langle a_1, \dots, a_n \rangle$  or simply as  $a_0$  if  $n = 0$ .

Let  $X = \{x_1, x_2, x_3, \dots\}$  be a set of **variables**. For  $K > 0$ , we shall often use the subsets  $X_K = \{x_1, \dots, x_K\}$ . A **ranked alphabet**  $\Sigma$  is a (potentially infinite) set disjoint from  $X$  paired with a function  $ar_\Sigma : \Sigma \rightarrow \mathbb{N}$  (or just  $ar$  where  $\Sigma$  is clear from context). The set of **ranked trees over  $\Sigma$  with variables in  $X$** , denoted by  $\mathcal{T}_\Sigma(X)$ , is the least set that contains  $X$  and all  $a \langle t_1, \dots, t_n \rangle$  where  $a \in \Sigma$ ,  $n = ar_\Sigma(a)$ ,  $t_1, \dots, t_n \in \mathcal{T}_\Sigma(X)$ .  $\mathcal{T}_\Sigma(\emptyset)$  is the set of **ground  $\Sigma$ -trees**, also written  $\mathcal{T}_\Sigma$ . Notions of position, substitution, etc., are all defined as usual.

We next recall the definitions of deterministic top-down tree automata ( $\text{DTTA}$ ) and deterministic top-down tree transducers ( $\text{DTOP}$ ).

**Definition 1.** A **quasi  $\text{DTTA}$**  is a tuple  $A = (\Sigma, Q, q_{\text{ini}}, rhs)$  such that  $\Sigma$  is a ranked alphabet,  $Q$  a finite set,  $q_{\text{ini}} \in Q$ , and  $rhs$  is a partial function that maps pairs  $(q, a) \in Q \times \Sigma$  to tuples of  $Q^{ar(a)}$ . A  $\text{DTTA}$  is a quasi  $\text{DTTA}$  for which  $\Sigma$  is finite, and thus so is  $rhs$ .

The elements of  $Q$  are called the states of  $A$ ,  $q_{\text{ini}}$  its initial state. The rules of  $A$  have the form  $q \langle a \langle x_1, \dots, x_n \rangle \rangle \rightarrow rhs(q, a)$ , where  $rhs(q, a)$  is defined and  $n = ar(a)$ . Each state  $q$  of a quasi  $\text{DTTA}$   $A$  recognises a tree language  $\llbracket q \rrbracket_A \subseteq \mathcal{T}(\Sigma)$ , (or just  $\llbracket q \rrbracket$  when  $A$  is clear from the context) defined by induction on the trees: we have  $a \langle t_1, \dots, t_n \rangle \in \llbracket q \rrbracket$  iff there is a rule  $q \langle a \langle x_1, \dots, x_n \rangle \rangle \rightarrow q_1, \dots, q_n$  and  $t_k \in \llbracket q_k \rrbracket$  for all  $k \in 1..n$ . The semantics of the automaton is  $\llbracket A \rrbracket = \llbracket q_{\text{ini}} \rrbracket$ .

Bottom-up tree automata ( $\text{BUTA}$ ) are defined similarly and as usual, with rules of the form  $a \langle q_1, \dots, q_n \rangle \rightarrow q$ .

**Definition 2.** A **quasi  $\text{DTOP}$**  is a tuple  $M = (\Sigma, \Delta, Q, ax, rhs)$  such that  $\Sigma$  and  $\Delta$  are ranked alphabets,  $Q$  is a finite set,  $q_{\text{ini}} \in Q$ , and  $rhs$  is a partial function that maps pairs  $(q, a) \in Q \times \Sigma$  to  $\mathcal{T}_\Delta(Q \times X_{ar(a)})$ . A  $\text{DTOP}$  is a quasi  $\text{DTOP}$  for which  $\Sigma$  and  $\Delta$  are finite, and thus  $rhs$  as well.

$\Sigma$  and  $\Delta$  are called input and output alphabets; the other components are as in automata. Each state  $q \in Q$  has as semantics a partial function  $\llbracket q \rrbracket$  from  $\mathcal{T}_\Sigma$  to  $\mathcal{T}_\Delta$ , defined by induction on terms  $t = a \langle t_1, \dots, t_n \rangle \in \mathcal{T}_\Sigma$  such that:

$$\llbracket q \rrbracket(t) = rhs(q, a) [q' \langle x_k \rangle \leftarrow \llbracket q' \rrbracket(t_k) \mid q' \in Q, k \in 1..ar(a)]. \quad (5)$$

The transformation defined by  $M$  is the partial function  $\llbracket M \rrbracket = \llbracket q_{\text{ini}} \rrbracket$ .

## 3 Symbolic Tree Automata and Transducers

In this section, we recall the definitions of symbolic  $\text{DTTA}$  and symbolic  $\text{DTOP}$  as in [16].

Symbolic machines are finite representations of potentially infinite quasi  $\text{DTTA}$  and quasi  $\text{DTOP}$ . They use **descriptors** to stand for the potentially infinite sets and functions. Given a set  $S$ , we call a set  $D$  paired with a function  $\llbracket . \rrbracket : D \rightarrow S$  as set of descriptors of elements of  $S$ . For instance, we can use the set  $E$  of regular expressions  $e \in E$  over an alphabet  $A$  as descriptors of regular languages  $\llbracket e \rrbracket \subseteq A^*$ . Outside of the definitions, we shall often assimilate the descriptors and their semantics.

**Definition 3.** A **symbolic  $\text{DTTA}$**  is a tuple  $A = (\Sigma, \Phi, Q, q_{\text{ini}}, rhs)$  such that  $(\Phi, Q, q_{\text{ini}}, rhs)$  is a quasi  $\text{DTTA}$  with a finite set of rules,  $\Phi$  is an alphabet of descriptors for subsets of the alphabet  $\Sigma$ , with  $ar(a) = ar(\varphi)$  for any  $a \in \llbracket \varphi \rrbracket$  and  $\varphi \in \Phi$ . For all  $a \in \Sigma$  and  $q \in Q$ , there exists at most one  $\varphi \in \Phi$  such that  $rhs(q, \varphi)$  is defined and  $a \in \llbracket \varphi \rrbracket$ .

The elements of  $\varphi \in \Phi$  are called (descriptors for) guards. A symbolic DTTA  $A$  is a finite representation of a (potentially) infinite quasi DTTA  $A'$  such that for every rule  $r$  of form,  $q\langle\varphi\langle x_1, \dots, x_n \rangle\rangle \rightarrow q_1, \dots, q_n$  of  $A$ , and for every  $\alpha \in \llbracket\varphi\rrbracket$ , there is a rule  $q\langle\alpha\langle x_1, \dots, x_n \rangle\rangle \rightarrow q_1, \dots, q_n$  in  $A'$ . The semantics of  $A$  is defined as that of  $A'$ : for all  $q \in Q$ ,  $\llbracket q \rrbracket_A = \llbracket q \rrbracket_{A'}$ . Symbolic BUTA are defined similarly.

**Definition 4.** A symbolic DTTA is **effective** if it satisfies the following conditions, which we always assume: **(1)** The set of guards  $\Phi$  is closed under conjunction and negation, i.e., there exists an algorithm computing some function  $\wedge : \Phi \times \Phi \rightarrow \Phi$  such that  $\llbracket \varphi \wedge \varphi' \rrbracket = \llbracket \varphi \rrbracket \cap \llbracket \varphi' \rrbracket$  for all  $\varphi, \varphi' \in \Phi$ , and an algorithm computing some function  $\neg : \Phi \rightarrow \Phi$  such that  $\llbracket \neg \varphi \rrbracket = \Sigma \setminus \llbracket \varphi \rrbracket$ . **(2)** There exists an algorithm deciding membership  $\alpha \in \llbracket \varphi \rrbracket$  given a guard  $\varphi \in \Phi$  and a label  $\alpha \in \Sigma$ ,

**Definition 5.** A **symbolic DTop** is a tuple  $M = (\Sigma, \Delta, \Phi, \mathcal{F}, Q, q_{\text{ini}}, \text{rhs})$  such that  $(\Phi, \mathcal{F}, Q, q_{\text{ini}}, \text{rhs})$  is a quasi DTop with a finite set of rules,  $\mathcal{F}$  is an alphabet of descriptors for partial functions from the input alphabet  $\Sigma$  to the output alphabet  $\Delta$ , with  $ar(\llbracket f \rrbracket(\alpha)) = ar(f)$  for every  $\alpha \in \text{dom}(f)$  and  $f \in \mathcal{F}$ . The same conditions on  $\Sigma, \Phi$  and  $\text{rhs}$  apply as for symbolic DTTA above.

The elements of  $f \in \mathcal{F}$  are called (descriptors for) data transformations. A symbolic DTop  $M$  is a finite representation of a (potentially) infinite quasi DTop  $M' = (\Sigma, \Delta, Q, \text{ax}, \text{rhs}')$ , such that for every rule  $q\langle\varphi\langle x_1, \dots, x_n \rangle\rangle \rightarrow \text{rhs}(q, \varphi)$  of  $M$ , and for every  $\alpha \in \llbracket\varphi\rrbracket$ , there is a rule  $q\langle\alpha\langle x_1, \dots, x_n \rangle\rangle \rightarrow \text{rhs}(q, \varphi) \mid f \leftarrow \llbracket f \rrbracket(\alpha) \mid f \in \mathcal{F}$  in  $M'$ . The semantics of  $M$  is defined as that of  $M'$ : for all  $q \in Q$ ,  $\llbracket q \rrbracket_M = \llbracket q \rrbracket_{M'}$ .

**Definition 6.** A symbolic DTop is **effective** (which is assumed in the remainder) if the underlying DTTA is, and **(1)** There is an algorithm that computes the value of the data transformation  $\llbracket f \rrbracket(\alpha)$  for a given  $f \in \mathcal{F}$  and  $\alpha \in \Sigma$  and returns  $\perp$  if it is not defined. **(2)** There is an algorithm that decides whether the image of a data transformation  $\llbracket f \rrbracket(\Sigma)$  is empty for a given  $f \in \mathcal{F}$ .

In **symbolic DTop<sup>R</sup>**, a symbolic BUTA on the transducer's input signature  $\Sigma$  first annotates the tree with its states  $P$ , and then the symbolic DTop transforms the annotated tree on  $\Sigma \times P$ .

Consider the logs of an application on a Unix-flavoured system. The log is a text file named "log", here containing "s" for every successful login, and "f" for every failure. Every week, the old log is discarded and replaced by statistics: the number of successful and failed logins (Figure 1). For this, we denote by  $f_{\text{stats}}$  the function counting the numbers  $n, m$  of occurrences of "s" and "f" in a string (here a log's contents), and outputting the string  $n";m$ . A fresh "log" is then created. The older log's statistics are named "stats.1", "stats.2", etc. so that higher numbers indicate older stats. To model this with a symbolic DTop<sup>R</sup>, we represent the contents of the logs folder by a list (CONS and NIL being the usual constructors) of files, each file being a tree of the form

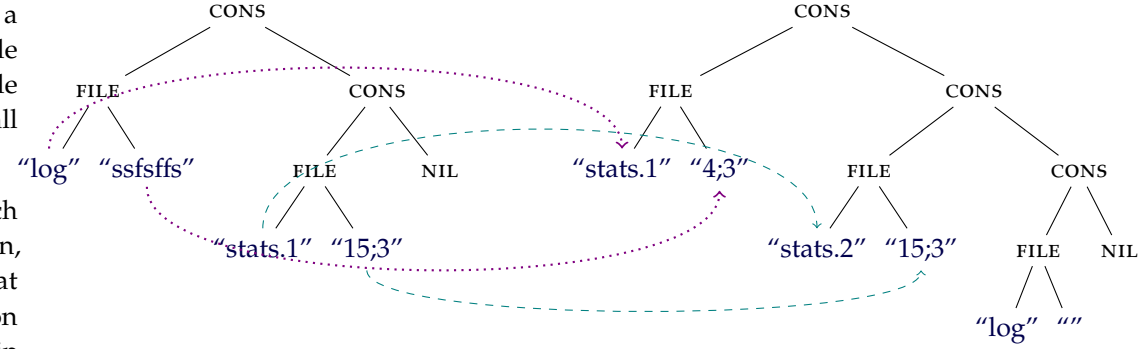


Fig. 1. Log cleanup and statistics: input tree on the left, output on the right.

FILE("filename", "contents"). The input and output alphabets are thus strings, along with FILE, CONS and NIL. The guards will be a small subset of regular expressions on strings plus descriptors matching CONS and NIL. The lookahead's purpose is to check whether the filename matches a stat file or not (which cannot be done in a top-down transducer's rule, in contrast to the extended rules (3) and (4)) and to annotate the FILE nodes with its findings. Guards are regular expressions. The descriptor matching a specific string is the string itself; \* matches everything; "stats."("0".."9")<sup>+</sup> is the descriptor matching stats filenames. The lookahead (LA) rules are

$$\begin{aligned}
 \text{"stats."("0".."9")}^+ \langle \rangle &\rightarrow p_{\text{stats}} & \text{FILE}\langle p_{\text{stats}}, p \rangle &\rightarrow p_{\text{stats}} \text{ file} \\
 \text{"log"} \langle \rangle &\rightarrow p & \text{FILE}\langle p, p \rangle &\rightarrow p \\
 \text{"s" | "f"}^* \langle \rangle &\rightarrow p & \text{CONS}\langle \_, p \rangle &\rightarrow p \quad \text{NIL}\langle \rangle \rightarrow p
 \end{aligned}$$

The label functions  $\mathcal{F}$  are taken as, for instance, the class of rational functions, which we can implement with word transducers with lookahead [6], satisfying all requisite properties. We represent a constant function by the string it produces, the identity by  $f_{\text{id}}$ , and  $f_{\text{increment}}$  for the function taking strings of the form "stats."k, where k is the decimal representation of an integer, and yielding "stats."(k + 1). We start in state  $q$ :

$$q\langle \text{NIL} : p \rangle \rightarrow \text{CONS}\langle \text{FILE}\langle \text{"log"}, \text{""} \rangle, \text{NIL} \rangle \quad (1')$$

$$q\langle \text{CONS} : p\langle x_1, x_2 \rangle \rangle \rightarrow \text{CONS}\langle q_{\text{name}}\langle x_1 \rangle, q\langle x_2 \rangle \rangle \quad (2')$$

$$q_{\text{name}}\langle \text{FILE} : p\langle x_1, x_2 \rangle \rangle \rightarrow \text{FILE}\langle q_{\text{log}}\langle x_1 \rangle, q_{\text{stats}}\langle x_2 \rangle \rangle \quad (3a)$$

$$q_{\text{name}}\langle \text{FILE} : p_{\text{stats file}}\langle x_1, x_2 \rangle \rangle \rightarrow \text{FILE}\langle q_{\text{incr}}\langle x_1 \rangle, q_{\text{id}}\langle x_2 \rangle \rangle \quad (4a)$$

$$q_{\text{incr}}\langle \text{"stats."("0".."9")}^+ : p_{\text{stats}} \rangle \rightarrow f_{\text{increment}} \quad (3b)$$

$$q_{\text{id}}\langle * : p \rangle \rightarrow f_{\text{id}} \quad q_{\text{stats}}\langle * : p \rangle \rightarrow f_{\text{stats}} \quad q_{\text{log}}\langle \text{"log"} : p \rangle \rightarrow \text{"stats.1"} \quad (4b)$$

## 4 Domain and Composition

To study problems like computation or equivalence on symbolic  $D_{\text{TOP}}$ , it is worth considering those problems as extensions of their counterparts for  $D_{\text{TOP}}$ . Indeed, most difficulties that previous papers [21, 20, 16] encountered are already relevant for the composition, normalization, or equivalence problems in the finite-labelled case [12, 17]. Most of those difficulties come from dealing with the domain of a transducer’s transformation. Since these problems have been solved in  $D_{\text{TOP}}$  and proofs in symbolic  $D_{\text{TOP}}$  are essentially identical, we shall only present the results, a reference for the proofs in  $D_{\text{TOP}}$ , and the additional conditions required of  $\Phi$  and  $\mathcal{F}$  for them to carry over to the symbolic case.

The first important results concern automata and their expressive power. Symbolic  $D_{\text{TTA}}$  which, as said before, we always assume to be effective, and have the classical properties of  $D_{\text{TTA}}$  (e.g. in [5]).

**Lemma 7.** (1) *The class of languages described by symbolic  $D_{\text{TTA}}$  is closed under intersection.*  
 (2) *If equivalence is decidable on  $\Phi$ , then equivalence is decidable on symbolic  $D_{\text{TTA}}$ .*

The second important result concerns the domains of symbolic  $D_{\text{TOP}}$ .

**Lemma 8.** *Let  $M$  be a symbolic  $D_{\text{TOP}}$ . Then we can build a symbolic  $D_{\text{TTA}}$   $A$  such that  $\llbracket A \rrbracket = \text{dom}(\llbracket M \rrbracket)$ .*

*Proof.* This is directly adaptable from the classical construction on  $D_{\text{TOP}}$ : see [8, 12] and [2, Lem. 39 p57]. The nub of the construction is that, as subtrees can be copied, each state of  $A$  corresponds to a set of states that explore the same subtree during the run. In the worst case scenario, this results in a full powerset construction, and thus leads to an exponential blowup in the number of states. Note that in the symbolic case, we need to compute the intersections of all the guards that may be invoked on the same subtree.  $\square$

The third important result concern the composition of symbolic  $D_{\text{TOP}}$ . In [16], the authors state that “*the proof of the first statement [of [21], that their algorithm computes the composition of any two symbolic  $D_{\text{TOP}}$ , and that the result is a symbolic  $D_{\text{TOP}}$ ] is insufficient.* [They proceed to show why that specific algorithm is incorrect] *We even conjecture that this statement is wrong, i.e., symbolic  $D_{\text{TOP}}$  are not closed under composition.*” Indeed they are not, and for the same reasons that  $D_{\text{TOP}}$  themselves are not closed under composition. This is mainly due to the problem of domain restriction.

**Example 9.** For  $a, b \in \mathcal{U}$ , it is easy enough to make a  $D_{\text{TOP}}$  that computes  $\tau_1$ , the identity over trees of form  $a(b, t_2)$ , and a  $D_{\text{TOP}}$  that computes  $\tau_2$ , that sends all trees of the form  $a(t_1, t_2)$  to their right subtree  $t_2$ . However, no  $D_{\text{TOP}}$  can compute  $\tau_2 \circ \tau_1$ , that sends all trees of form  $a(b, t_2)$  to their right subtree  $t_2$ : such a transducer would

need to start in a state  $q_0$  that only has one rule of form  $q_0 \langle a \langle x_1, x_2 \rangle \rangle \rightarrow t$ . This rule cannot produce any output symbol, and thus  $t$  is of the form  $q \langle x_i \rangle$ . This  $x_i$  must be  $x_1$ , as the transducer needs to visit the left subtree to ensure that it is  $b$ . This  $x_i$  must also be  $x_2$ , as the transducer needs to copy the right subtree. Hence, no such  $D_{\text{TOP}}$  can exist.

Note that this argument applies equally to symbolic  $D_{\text{TOP}}$  and to  $D_{\text{TOP}}$ ; only the general form of rules matter. To find a class closed under composition, a solution presented in the  $D_{\text{TOP}}$  case [12] is to consider transducers *with domain inspection*: a **symbolic  $D_{\text{TOP}}$  with inspection** is a pair  $N = (M, A)$  of a symbolic  $D_{\text{TOP}}$   $M$  and a symbolic  $D_{\text{TTA}}$   $A$ . Its semantic is  $\llbracket N \rrbracket = \llbracket M \rrbracket_{\llbracket A \rrbracket}$ , the function of  $M$  restricted to the language of  $A$ . We know that  $D_{\text{TOP}}$  with inspection are closed under composition [12]. This result extends to symbolic  $D_{\text{TOP}}$  if the set of functions of  $\mathcal{F}$  is itself closed under composition, and the images of guards through functions of  $\mathcal{F}$  form a suitable set of guards (i.e. they satisfy the requirements for effectiveness).

**Lemma 10.** *Let  $\mathcal{F}$  be closed under composition, and  $N, N'$  be two effective symbolic  $D_{\text{TOP}}$  with inspection using functions of  $\mathcal{F}$ . Then we can build a symbolic  $D_{\text{TOP}}$  with inspection  $N''$  such that  $\llbracket N'' \rrbracket = \llbracket N \rrbracket \circ \llbracket N' \rrbracket$ .*

The main intuition behind the generalisation of the classical results [1] is presented in several papers [21, 20, 16]; roughly, a rule in  $N''$  is the image of the right-hand side of a rule of  $N'$  by a state of  $N$ . To obtain closure by composition for symbolic  $D_{\text{TOP}}$  – or indeed for  $D_{\text{TOP}}$ , as the problem is fundamentally unchanged by the alphabets – necessitates the use of either very strong restrictions, as in [16], or the use of domain inspection, which we prefer here.

## 5 Deciding Equivalence

In this section, we show that, given a few basic properties on label transformations (mostly that equivalence is decidable for label transformations) the equivalence problem for symbolic  $D_{\text{TOP}}$  is decidable, regardless of linearity, by reducing that problem to equivalence for  $D_{\text{TOP}}$ , which is known to be decidable.

The main observation behind the reduction is that, although  $\Sigma$ ,  $\Phi$  and  $\mathcal{F}$  may well be infinite, only a finite number of predicates and transformations are actually used in a symbolic  $D_{\text{TOP}}$  (or indeed, any pair of symbolic  $D_{\text{TOP}}$ ), and give rise to a finite number of behaviours. These finite subsets of  $\Phi$  and  $\mathcal{F}$  will serve as finite input and output signatures in our reduction, which encodes those finitely many behaviours.

We shall first present a straightforward but coarse reduction, which captures *most* of the ways in which two symbolic  $D_{\text{TOP}}$  may differ, with the exception of some more technical cases where the *origins* of produced nodes matter. The notion of origin will then be presented and used to complete the construction.

## 5.1 The Basic Reduction

The first reduction is as follows: let  $M = (\Phi, \Sigma, \mathcal{F}, \Delta, P, p_{\text{ini}}, R)$  and  $N = (\Phi, \Sigma, \mathcal{F}, \Delta, Q, q_{\text{ini}}, S)$  be two symbolic DTop. We build their DTop **representations**, the DTop  $\underline{M}$  and  $\underline{N}$ . Strictly speaking we should write  $\underline{M}^{M,N}$  and  $\underline{N}^{M,N}$ , as the construction is specific to the pair of transducers under consideration, and the same applies to the representation of each component of the transducers –  $\underline{\Phi}$ ,  $\underline{\Sigma}$ , etc – which we define below. In this section we assimilate descriptors  $\varphi, f$  and their semantics  $\llbracket \varphi \rrbracket, \llbracket f \rrbracket$  to lighten the notations.

We make the following additional **equivalence-testing assumptions**: for all  $\varphi, \psi \in \Phi$ , it is decidable whether  $\varphi = \psi$ . For all  $f, g \in \mathcal{F}$  and all  $\varphi \in \Phi$ , it is decidable whether there exists some  $c \in \Delta$  such that  $f(\varphi) = \{c\}$ , and this  $c$  is computable; and it is decidable whether  $f|_{\varphi} = g|_{\varphi}$ .

The finite information relevant to the behaviour of symbolic DTop is which guards are satisfied. Thus we let

$$\Pi = \{\text{gd}(r) \mid r \in R \cup S\} \subseteq \Phi \quad (6)$$

be the subsets of guards actually used by either one of the two transducers.

The **representation of  $a \in \Sigma$**  is the equivalence class of  $a \in \Sigma$  with respect to the tests of  $\Pi$ ; that is to say, the set of all labels which are indistinguishable from  $a$  by our transducers:

$$\underline{a} = \{b \in \Sigma \mid \forall \pi \in \Pi, a \in \pi \Leftrightarrow b \in \pi\}, \quad (7)$$

which is computed as

$$\underline{a} = \bigcap_{\pi \in \Pi} \pi \setminus \bigcup_{a \notin \pi} \pi. \quad (8)$$

The **finite alphabet  $\underline{\Sigma}$  representing  $\Sigma$**  is defined as  $\underline{\Sigma} = \{\underline{a} \mid a \in \Sigma\}$ . Note that this is indeed a finite set, as  $\Pi$  is finite.

To properly represent a guard  $\varphi \in \Pi$ , we must make sure to capture all the possible circumstances in which the guard is used; recall from Lemma 8<sub>[p4]</sub> that a subtree may be evaluated in any number of states during a run, and thus simultaneously tested by any number of guards, each combination corresponding to a different behaviour of the transducer, and to a new, finer guard. Thus the representation of a guard  $\varphi \in \Pi$  must be the set of all the (finer) guards obtained by such combinations. Those combinations are precisely the equivalence classes of all labels that satisfy  $\varphi$ . Thus the **representation of a guard  $\varphi \in \Pi$**  is defined as the set of equivalence classes of labels accepted by  $\varphi$ :

$$\underline{\varphi} = \{\underline{a} \mid a \in \Sigma, a \in \varphi\}. \quad (9)$$

The **representation of an input tree  $t \in \mathcal{T}(\Sigma)$**  is defined inductively as

$$\underline{a}(t_1, \dots, t_n) = \underline{a}(\underline{t}_1, \dots, \underline{t}_n) \in \mathcal{T}(\underline{\Sigma}). \quad (10)$$

For the representation of a label transformation  $f$ , it is important to store the *actual* transformation that can be performed during a run, which means restricting by the appropriate guard. There is also a more subtle technicality. We usually store the function itself, except in the case where only a constant can be produced: then one must store only the constant itself, regardless of the input domain. This is linked to the question of the *origin* of a node, which will be explained in the next sections. Thus the **representation of a label transformation  $f$  restricted to  $\varphi$**  is defined as

$$\underline{f}|_{\varphi} = \begin{cases} c & \text{if } f(\varphi) = \{c\}, \text{ and} \\ f|_{\varphi} & \text{otherwise.} \end{cases} \quad (11)$$

The **representation of a rule  $r \in R \cup S$** , of the form  $r = q\langle\varphi\langle x_1, \dots, x_n \rangle\rangle \rightarrow t$ , is given by the set  $\underline{r}$  of all classical rules

$$q\langle\rho\langle x_1, \dots, x_n \rangle\rangle \rightarrow t[f \leftarrow \underline{f}|_{\rho} \mid f \in \mathcal{F}], \quad (12)$$

where  $\rho \in \underline{\varphi}$ . Letting  $\underline{R} = \bigcup_{r \in R} \underline{r}$  and  $\underline{S} = \bigcup_{s \in S} \underline{s}$ , we finally have  $\underline{M} = (\underline{\Sigma}, \underline{\Delta}, P, p_{\text{ini}}, \underline{R})$  and  $\underline{N} = (\underline{\Sigma}, \underline{\Delta}, Q, q_{\text{ini}}, \underline{S})$ .

## 5.2 Examples; The Necessity of Origins

It is easy to see that if  $\underline{M}$  and  $\underline{N}$  are not equivalent, then neither are  $M$  and  $N$ , as a difference in the representation means that the symbolic DTop must either take in different inputs, produce tree with different shapes, or produce different labels. Let us visualise that on an example:

**Example 11.** Consider  $M$  of initial state  $p$  and rules

$$p\langle\text{"c"}^*\langle x_1, x_2 \rangle\rangle \rightarrow f_{\text{id}}\langle p\langle x_1 \rangle, p\langle x_2 \rangle \rangle \quad (13)$$

$$p\langle\text{"a"}^*\rangle \rightarrow f_{\text{id}} \quad p\langle\text{"b"}^*\rangle \rightarrow f_{\text{id}} \quad (14)$$

and  $N$  of initial state  $q$  and rules

$$q\langle\text{"c"}^*\langle x_1, x_2 \rangle\rangle \rightarrow f_{\text{id}}\langle q\langle x_2 \rangle, q\langle x_1 \rangle \rangle \quad (15)$$

$$q\langle\text{"a"}^*\rangle \rightarrow f_{\text{id}} \quad q\langle\text{"b"}^*\rangle \rightarrow f_{\text{id}} \quad (16)$$

The representation of  $M$  is easily computed, as the guards of  $\Pi$  are disjoint, and describes the transformation

$$\llbracket \underline{M} \rrbracket = \begin{array}{c} \text{"c"}^* \\ / \quad \backslash \\ \text{"a"}^* \quad \text{"b"}^* \end{array} \mapsto \begin{array}{c} f_{\text{id}} | \text{"c"}^* \\ / \quad \backslash \\ f_{\text{id}} | \text{"a"}^* \quad f_{\text{id}} | \text{"b"}^* \end{array}. \quad (17)$$

In that case, this is even a finite transformation – there is a single possible input tree. Likewise for  $\mathbb{N}$ :

$$\llbracket \mathbb{N} \rrbracket = \begin{array}{c} \text{"c"*} \\ / \quad \backslash \\ \text{"a"*} \quad \text{"b"*} \end{array} \mapsto \begin{array}{c} f_{id}|\text{"c"*} \\ / \quad \backslash \\ f_{id}|\text{"b"*} \quad f_{id}|\text{"a"*} \end{array} . \quad (18)$$

We have  $\llbracket \mathbb{M} \rrbracket \neq \llbracket \mathbb{N} \rrbracket$ , as  $f_{id}|\text{"a"*} \neq f_{id}|\text{"b"*}$ , and this reflects the difference of behaviour between  $\mathbb{M}$  and  $\mathbb{N}$ .

**Example 12.** Consider now  $\mathbb{M}$  of initial state  $p$  and rules

$$\begin{aligned} p\langle \text{"c"*}(x_1, x_2) \rangle &\rightarrow f_{id}\langle p\langle x_1 \rangle, p\langle x_2 \rangle \rangle \\ p\langle \text{"a"*} \rangle &\rightarrow \text{"a"} \quad p\langle \text{"b"*} \rangle \rightarrow \text{"a"} \end{aligned} \quad (19)$$

and  $\mathbb{N}$  of initial state  $q$  and rules

$$\begin{aligned} q\langle \text{"c"*}(x_1, x_2) \rangle &\rightarrow f_{id}\langle q\langle x_2 \rangle, q\langle x_1 \rangle \rangle \\ q\langle \text{"a"*} \rangle &\rightarrow \text{"a"} \quad q\langle \text{"b"*} \rangle \rightarrow \text{"a"} \end{aligned} \quad (20)$$

Since we produce constants, by (11)<sub>[p5]</sub> we store this time the constant itself, and not the restricted function; thus we have the representations:

$$\llbracket \mathbb{M} \rrbracket = \llbracket \mathbb{N} \rrbracket = \begin{array}{c} \text{"c"*} \\ / \quad \backslash \\ \text{"a"*} \quad \text{"b"*} \end{array} \mapsto \begin{array}{c} f_{id}|\text{"c"*} \\ / \quad \backslash \\ \text{"a"} \quad \text{"a"} \end{array} . \quad (21)$$

This is good, as in that case, it does not matter whether the constant "a" was obtained from an input label in "a"\* or in "a":  $\mathbb{M}$  and  $\mathbb{N}$  are equivalent regardless. Had we stored the constant function with its domain restrictions instead of the constant itself, we would have obtained  $\llbracket \mathbb{M} \rrbracket \neq \llbracket \mathbb{N} \rrbracket$ , and thus a false negative. However, there is work left to do in order to obtain the reciprocal property; that is to say, to ensure that if the representations are equivalent, then so must the originals. A variation of Example 11 provides the counter-example:

**Example 13.** Consider  $\mathbb{M}$  of initial state  $p$  and rules

$$\begin{aligned} p\langle \text{"c"*}(x_1, x_2) \rangle &\rightarrow f_{id}\langle p\langle x_1 \rangle, p\langle x_2 \rangle \rangle \\ p\langle \text{"a"*} \rangle &\rightarrow f_{id} \quad p\langle \text{"a"*} \rangle \rightarrow f_{id} \end{aligned} \quad (22)$$

and  $\mathbb{N}$  of initial state  $q$  and rules

$$\begin{aligned} q\langle \text{"c"*}(x_1, x_2) \rangle &\rightarrow f_{id}\langle q\langle x_2 \rangle, q\langle x_1 \rangle \rangle \\ q\langle \text{"a"*} \rangle &\rightarrow f_{id} \quad q\langle \text{"a"*} \rangle \rightarrow f_{id} \end{aligned} \quad (23)$$

We have:

$$\llbracket \mathbb{M} \rrbracket = \llbracket \mathbb{N} \rrbracket = \begin{array}{c} \text{"c"*} \\ / \quad \backslash \\ \text{"a"*} \quad \text{"a"*} \end{array} \mapsto \begin{array}{c} f_{id}|\text{"c"*} \\ / \quad \backslash \\ f_{id}|\text{"a"*} \quad f_{id}|\text{"a"*} \end{array} . \quad (24)$$

Yet it is not the case that  $\mathbb{M}$  and  $\mathbb{N}$  are equivalent: consider the input tree  $t = \text{"c"}\langle \text{"a"}, \text{"aa"} \rangle$ : we have  $\llbracket \mathbb{M} \rrbracket(t) = t$ , whereas  $\llbracket \mathbb{N} \rrbracket(t) = \text{"c"}\langle \text{"aa"}, \text{"a"} \rangle$ . Though the same operation is done on the two children, that does not make them interchangeable. Thus we must, additionally, have in the representation some means of tracking the *origin* is output labels; that is to say, from which node they are produced – left child or right child, in that case. This origin information must apply even if the domains of the transformations are identical.

Yet it must be noted that, if we combined this example with the previous one, producing the same constant left and right, we would *again* be in a case where no origin information should be stored, as it won't matter whether the constant "comes from" the right or left child.

The next section formally introduces the notion of *origin*.

### 5.3 Origins and Origin Equivalence

The problem identified in Example 13 is that a reduction from symbolic  $\text{DTOP}$  to  $\text{DTOP}$  can only be accurate if it stores not only the information relative to which guard or function was used to transform input labels into output labels, but also enough information to know which input label is transformed into which output label. This notion exists in previous papers: notably, [3] discusses *origin*, a relation between the input and output nodes of a transformation. We will define similar notions, identifying for each output node the input node, guard, and function used for its production.

We define a position to be a word on the alphabet  $\mathbb{N}$ . We can identify each node of a tree  $s \in \mathcal{T}_{\Sigma}$  with a unique position.

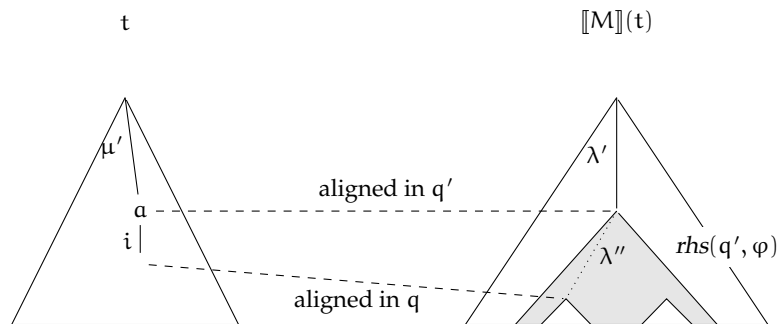
- ◊ The root of  $s$  is at the position  $\varepsilon$ .
- ◊ If a node of  $s$  is at position  $\mu$  and labeled by  $a \in \Sigma$  of arity  $k$ , then its  $i^{\text{th}}$  child has the position  $\mu i$ , where  $1 \leq i \leq k$ .

For instance, in Figure 1, the node 1 is labelled by `FILE`, and the node 12 in the input is labelled by `"ssfsffs"`.

We use this notion of position to talk of *syntactic alignment*. This notion, found in [2], aims to describe which part of the input is read by a symbolic  $\text{DTOP}$  to produce a given part of the output. It is, as its name implies, a purely syntactic property, describing the functioning of the transducer. We define this notion recursively:

**Definition 14.** Let  $M$  be a symbolic  $\text{DTOP}$ , and  $t \in \text{dom}(\llbracket M \rrbracket)$  a tree. We define by induction what it means for  $\mu$  a position of  $t$  and  $\lambda$  a position of  $\llbracket M \rrbracket(t)$  to be syntactically aligned in state  $q$  (we write  $\mu \sim_q \lambda$ ):

- ◇  $\mu = \varepsilon$  and  $\lambda = \varepsilon$  are aligned in state  $q_{\text{ini}}$ .
- ◇ If  $\mu' \sim_{q'} \lambda'$ , the node at position  $\mu'$  in  $t$  is labeled by  $a$ , and there exists a rule  $q' \langle \varphi \langle x_1, \dots, x_n \rangle \rangle \rightarrow \text{rhs}(q', \varphi)$  such that  $a \in \varphi$ , and furthermore there is a node of  $\text{rhs}(q', \varphi)$  at position  $\lambda''$  labeled by  $q \langle x_i \rangle$ , then  $\mu' i$  and  $\lambda' \lambda''$  are aligned in state  $q$ .



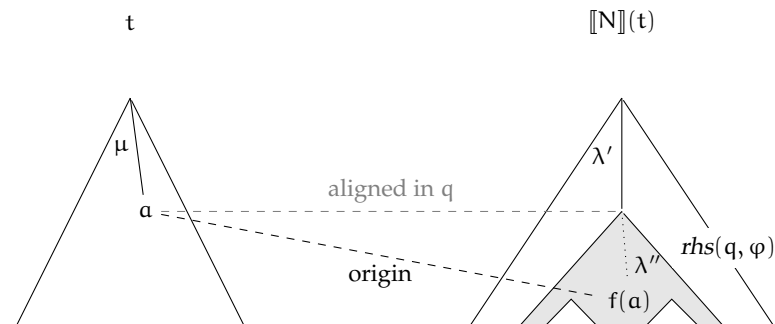
We use the notion of syntactic alignment to define a notion of origin similar to the one found in [3]. The intuition of this notion goes as follows: let  $M$  be a symbolic  $\text{DTOP}$ ,  $t \in \text{dom}(\llbracket M \rrbracket)$ , and  $\lambda$  a position in  $\llbracket M \rrbracket(t)$ . We want to know which node of  $t$  was read by  $M$  when the label at position  $\lambda$  was produced in  $\llbracket M \rrbracket(t)$ . We call this input node the origin of the output node at position  $\lambda$ . Note that in the symbolic case, we also are interested to track the guard of the rule used to produce said node, and the function used to transform an input label in the output label at position  $\lambda$ .

**Definition 15.** Let  $M$  be a symbolic  $\text{DTOP}$ ,  $t \in \text{dom}(\llbracket M \rrbracket)$  a tree, and  $\lambda$  a position of  $\llbracket M \rrbracket(t)$ . We say that  $\mu$  is the origin of  $\lambda$  if:

- ◇  $\mu$  is the position of a node of  $t$  labeled by some  $a \in \Sigma$
- ◇ there exists  $\lambda', \lambda''$  such that  $\lambda = \lambda' \lambda''$  and  $\mu \sim_q \lambda'$
- ◇ there exists a rule  $q \langle \varphi \langle x_1, \dots, x_n \rangle \rangle \rightarrow \text{rhs}(q, \varphi)$  such that  $a \in \varphi$ , and  $\lambda''$  is the position of a node of  $\text{rhs}(q, \varphi)$  labeled by some  $f \in \mathcal{F}$ .

We call  $f$  the *origin function* of  $\lambda$ , and  $\varphi$  the *origin guard* of  $\lambda$ .

Note that for a symbolic  $\text{DTOP}$   $M$  and a tree  $t$ , each position  $\lambda$  of  $\llbracket M \rrbracket(t)$  has exactly one origin: each output node is produced during one and only one step of the computation of  $\llbracket M \rrbracket(t)$ . For instance, in Example 13, the node at position 1 in an output of  $M$  comes from the input node at position 1. Its origin guard is “ $a$ ” $*$  and its origin function is



$f_{\text{id}}|_{\text{“}a\text{”}^*}$ . Conversely, in  $N$ , the node at position 1 in an output of  $M$  also has the origin guard “ $a$ ” $*$  and the origin function  $f_{\text{id}}|_{\text{“}a\text{”}^*}$ , but its origin is the input node of position 2. This difference in origins is not detected in the reduction attempt of Section 5.1.

In [3], the notion of origin is accompanied by a notion of origin equivalence. In essence, two transducers are origin-equivalent if they are equivalent and every output node is produced from the same origins in both transducers.

**Definition 16.** Two equivalent symbolic  $\text{DTOP}$   $M, N$  are *origin-equivalent* iff for all  $t \in \text{dom}(\llbracket M \rrbracket)$ , for any node  $\lambda$  of  $\llbracket M \rrbracket(t)$ , the origins of  $\lambda$  for  $M$  and  $N$  are identical.

The advantage of origin equivalence is that if the origins are known, origin equivalence becomes easier to test than equivalence (see [3]). Origin equivalence is stronger than equivalence: completing the reduction of 5.1 to embed information about the origin would prevent the false positives presented in Example 13, in which  $M$  and  $N$  are not equivalent. Although  $\underline{M}$  and  $\underline{N}$  are equivalent, they are not origin-equivalent: the output node at position 1 is of origin 1 for  $M$  and  $\underline{M}$ , but 2 for  $N$  and  $\underline{N}$ .

Unfortunately, origin equivalence is strictly stronger than equivalence, which means that by getting rid of false positives, we might introduce false negatives. In Example 12,  $M$  and  $N$  are equivalent, and, again, although  $\underline{M}$  and  $\underline{N}$  are equivalent, they are not origin-equivalent: the output node at position 1 is of origin 1 for  $M$  and  $\underline{M}$ , but 2 for  $N$  and  $\underline{N}$ .

To prevent this kind of false negative, we do not encode *all* origins, as we do not want to fully test for origin equivalence. Instead, we relax the notion of origin equivalence: the origin of an output node only matters if its origin function can produce different results for different input labels. If the origin function is constant, as in Example 12, then the origin node is of no relevance to determine equivalence. We define the notion of weak-origin equivalence as origin equivalence everywhere except when constants are produced.

**Definition 17.** Two equivalent symbolic  $\text{DTOP}$   $M, N$  are *weak-origin-equivalent* iff for all  $t \in \text{dom}(\llbracket M \rrbracket)$ , for any node  $\lambda$  of  $\llbracket M \rrbracket(t)$ , the origin nodes of  $\lambda$  for  $M$  and  $N$  are identical, or its origin functions for  $M$  and  $N$  are constant functions of same image.



Note that when we consider these origin functions, we consider them restricted to the origin guard. In the case of a rule  $q(\varphi\langle x_1, \dots, x_n \rangle) \rightarrow rhs(q, \varphi)$ , the function we want to test as constant is in fact  $f|_{\varphi}$ , not  $f$  itself. However, it is quite possible that equivalent transducers use different tests to compute the same transformations. This would needlessly affect the notion of weak-origin equivalence. To ensure no such problem arises, we introduce the notion of harmonized symbolic  $D_{TOP}$ , to say that they use the same guards to read the same nodes.

**Definition 18.** Two symbolic  $D_{TOP}$   $M, N$  are *harmonized* iff for all  $t \in \text{dom}[\llbracket M \rrbracket] \cap \text{dom}[\llbracket N \rrbracket]$ , for any node  $\lambda$  of  $\llbracket M \rrbracket(t)$  and  $\lambda'$  of  $\llbracket N \rrbracket(t)$ , if the origin of  $\lambda$  for  $M$  and  $\lambda'$  of  $N$  are identical, then their origin guards are identical.

It is easy to see that any two symbolic  $D_{TOP}$  can be harmonized. A way to do so would be to render all guards disjoint, similarly to (9)<sub>[p5]</sub> by computing all possible combinations of guards from the two transducers. While this method is not parsimonious, it shows that any pair of symbolic  $D_{TOP}$  can be harmonized without any change to their semantics by dividing their rules in small enough components.

We know that weak-origin equivalence is strictly weaker than origin equivalence. As seen previously in Example 12,  $\underline{M}$  and  $\underline{N}$  are not origin-equivalent as the output node at position 1 is of origin 1 for  $M$  and  $\underline{M}$ , but 2 for  $N$  and  $\underline{N}$ . However, this distinction disappears with weak-origin equivalence, as the origin function of 1 in both  $M$  and  $N$  is constant of image “a”:  $M$  and  $N$  are weak-origin-equivalent.

By definition, weak-origin equivalence is at least as strong as equivalence. We will show that for harmonized symbolic  $D_{TOP}$ , both these notions are, in fact, identical.

**Lemma 19.** *Two harmonized symbolic  $D_{TOP}$  are equivalent if and only if they are weak-origin-equivalent.*

*Proof.* Since weak-origin equivalence assumes equivalence, we only have to prove that two equivalent symbolic  $D_{TOP}$  are necessarily weak-origin-equivalent.

Let  $M, N$  be two equivalent symbolic  $D_{TOP}$ , and suppose they are not weak-origin-equivalent. There is a tree  $t \in \text{dom}(\llbracket M \rrbracket)$ , and a node  $\lambda$  of  $\llbracket M \rrbracket(t)$  such that  $\mu$  is its origin node for  $M$  and a different node  $\mu'$  is its origin node for  $N$ . Furthermore, one of its origin functions (we pick  $f$  its origin function for  $M$ , and  $\varphi$  its origin guard) is not constant. Since  $f$  is not constant we can replace the label at position  $\mu$  in  $t$  from  $a \in \varphi$  to another  $a' \in \varphi$  such that  $f(a) \neq f(a')$ . By replacing *only* the label of  $\mu$  in  $t$  we thus build another tree  $t'$ . Since  $M$  and  $N$  are harmonized,  $N$  uses the same guard  $\varphi$  as  $M$  to read the input label at  $\mu$ : thus this substitution bears no incidence on the rules used to compute  $\llbracket N \rrbracket(t')$ . The label of  $\lambda$  in  $\llbracket M \rrbracket(t')$  changed from  $f(a)$  to  $f(a')$ . However, since the label of  $\mu'$  has not changed in the input, the label of  $\lambda$  has not changed in  $\llbracket N \rrbracket(t')$ . This means  $\llbracket M \rrbracket(t') \neq \llbracket N \rrbracket(t')$ , which contradicts the equivalence of  $M$  and  $N$ .  $\square$

## 5.4 An Origin-Aware Reduction

We extend the reduction of 5.1 to include enough origin information to translate weak-origin equivalence. As we have seen previously, the origin of an output node is relevant if and only if a change of input label yields a change in output labels. To reduce this property to a finite signature, we create two distinct input labels for each guard  $\varphi$ , which will produce one of two distinct output labels for each non-constant function  $f$ . First, we create for each guard two distinct labels by placing, on each input node, a bit: the **representations of an input tree**  $t \in \mathcal{T}(\Sigma)$  becomes

$$\underline{a}(t_1, \dots, t_n) = \{ (\underline{a}, b)(u_1, \dots, u_n) \mid b \in \mathbb{B}, u_i \in \underline{t}_i, \forall i \} \subseteq \mathcal{T}(\Sigma \times \mathbb{B}), \quad (29)$$

which is as before with the addition of bits  $b$ , called **obit** (origin bit). Any of these bits can be set at 0 or 1, leading to several representations for each tree of  $\mathcal{T}(\Sigma)$ . To encode origin, this bit is carried over to the output by the rules of  $\underline{M}$ , but only when a change of input labels yields a change in output labels. Accordingly, the **representation of a label transformation  $f$  restricted to  $\varphi$ , with obit  $b$**  is redefined as

$$\underline{f}|_{\varphi, b} = \begin{cases} c & \text{if } f(\varphi) = \{c\}, \text{ and} \\ (f|_{\varphi}, b) & \text{otherwise,} \end{cases} \quad (30)$$

and the **representation of a rule  $r \in R \cup S$** , of the form  $r = q(\varphi\langle x_1, \dots, x_n \rangle) \rightarrow t$ , is given by the set  $\underline{r}$  of all classical rules

$$q(\langle \rho, b \rangle \langle x_1, \dots, x_n \rangle) \rightarrow t[f \leftarrow \underline{f}|_{\rho, b} \mid f \in \mathcal{F}], \quad (31)$$

for all  $b \in \mathbb{B}, \rho \in \varphi$ .

Consider the effect of this change on Example 13: we now have

$$\llbracket \underline{M} \rrbracket \ni \begin{array}{c} \text{“c”}^*, 0 \\ \swarrow \quad \searrow \\ \text{“a”}^*, 1 \quad \text{“a”}^*, 0 \end{array} \mapsto \begin{array}{c} f_{id} | \text{“c”}^*, 0 \\ \swarrow \quad \searrow \\ f_{id} | \text{“a”}^*, 1 \quad f_{id} | \text{“a”}^*, 0 \end{array}, \quad (32)$$

whereas

$$\llbracket \underline{N} \rrbracket \ni \begin{array}{c} \text{“c”}^*, 0 \\ \swarrow \quad \searrow \\ \text{“a”}^*, 1 \quad \text{“a”}^*, 0 \end{array} \mapsto \begin{array}{c} f_{id} | \text{“c”}^*, 0 \\ \swarrow \quad \searrow \\ f_{id} | \text{“a”}^*, 0 \quad f_{id} | \text{“a”}^*, 1 \end{array}, \quad (33)$$

and thus, properly,  $\llbracket \underline{M} \rrbracket \neq \llbracket \underline{N} \rrbracket$ .

A convenient way of thinking about obits is to envision the input trees where a single label bears a 1 as obit; then the set of output labels bearing a 1 as obit is exactly the set of nodes which have that input node as origin. For instance, in Figure 2, two output



**Fig. 2. Using obits to deduce origins.** We represent the obits with colours:  $\circ$  are nodes of obit 0,  $\bullet$  are nodes of obit 1. We apply two transformations  $\tau_1$  and  $\tau_2$  on an input tree (here in the middle):  $\tau_1$  replaces its right leaf by a copy of the left one, while  $\tau_2$  replaces its left leaf by a copy of the right one. Visually, “turning on” an obit in the input makes it possible to visualise all output nodes that depend upon it.

nodes share the same origin. Of course, we do not always store the origin information: in Example 12, we would still store only the constant, and discard information about both the transformation and the origin, thus avoiding a false negative.

This construction leads naturally to some interesting properties between  $\underline{M}$  and  $M$ , notably that the domains and origins are properly translated.

**Lemma 20.** *Let  $M$  be a symbolic  $\text{DTOP}$  and  $\underline{M}$  its representation.*

- $\diamond \text{dom}[\underline{M}] = \bigcup_{t \in \text{dom}[M]} \underline{t}$
- $\diamond$  for every  $s \in \underline{t}$ , if  $\mu$  is the origin of  $\lambda$  in  $\underline{M}$ , then  $\mu$  is the origin of  $\lambda$  in  $M$ .

The fact that  $\underline{M}$  translates the functioning of  $M$  in a finite signature means that we can now reduce the equivalence problem in symbolic  $\text{DTOP}$  to the equivalence problem in  $\text{DTOP}$ .

**Theorem 21.** *Let  $M, N$  be two symbolic  $\text{DTOP}$ . Then  $[\underline{M}] = [\underline{N}]$  if and only if  $[M] = [N]$ .*

*Proof.* First, by Lemma 20, the domain of  $\underline{M}$  is the set of all the representations of trees of  $\text{dom}([\underline{M}])$ . Hence  $\underline{M}$  and  $\underline{N}$  are of same domain if and only if  $M$  and  $N$  are of same domain.

Let us now assume that the domains of  $M$  and  $N$  are the same. We suppose they are harmonized, as described under Definition 18, by cutting all guards into the fragments described in (9)<sub>[p5]</sub>.

Suppose  $M$  and  $N$  are not equivalent; let  $t$  be an input tree such that  $u = [M](t) \neq [N](t) = v$ . We consider a position  $\lambda$  that exists both in  $u$  and  $v$  but where the label at  $\lambda$  differs in  $u$  and  $v$ . In  $M$  the origin node of  $\lambda$  is  $\mu$ ; the origin function is  $f|_\varphi$ . In  $N$  the origin node of  $\lambda$  is  $\mu'$ ; the origin function is  $g|_\psi$ .

If  $\mu = \mu'$  and  $f|_\varphi = g|_\psi$ , there would not be a difference between the labels of  $u$  and  $v$  at position  $\lambda$ . Hence of these equalities must not hold. Another point worth noting is that since the labels of  $u$  and  $v$  at position  $\lambda$  differ, then  $f|_\varphi$  and  $g|_\psi$  cannot be constant functions of same image.

If  $\mu \neq \mu'$ , i.e. if the origins of  $\lambda$  are different for  $M$  and  $N$ , then they also are for  $\underline{M}$  and  $\underline{N}$  for any  $s \in \underline{t}$  by Lemma 20. We choose an  $s$  whose label at  $\mu$  is  $(\varphi, 0)$  and

whose label at  $\mu'$  is  $(\psi, 1)$ . Then the label at  $\lambda$  in  $[\underline{M}](s)$  is  $f|_\varphi, 0$  while the label at  $\lambda$  in  $[\underline{N}](s)$  is  $g|_\psi, 1$ . The only way both of these labels could be the same despite the obit difference would be for  $f|_\varphi$  and  $g|_\psi$  to be constant of same image, which we know to be in contradiction with our suppositions. Hence in this case  $[\underline{M}] \neq [\underline{N}]$ .

If  $f|_\varphi \neq g|_\psi$ , i.e. if the origin functions of  $\lambda$  are different for  $M$  and  $N$ , then for any  $s \in \underline{t}$  the label at  $\lambda$  in  $[\underline{M}](s)$  is of form  $f|_\varphi, b$  with some  $b \in \{0, 1\}$ , while the label at  $\lambda$  in  $[\underline{N}](s)$  is of form  $g|_\psi, b'$  with some  $b' \in \{0, 1\}$ . The only way both these labels could be the same despite the function difference would be for  $f|_\varphi$  and  $g|_\psi$  to be constant functions of same image, which is again in contradiction with our suppositions. Hence in this case  $[\underline{M}] \neq [\underline{N}]$ .

Conversely, suppose  $\underline{M}$  and  $\underline{N}$  are not equivalent; let  $s$  be an input tree such that  $u' = [\underline{M}](s) \neq [\underline{N}](s) = v'$ . We consider a position  $\lambda$  that exists both in  $u'$  and  $v'$  but where the label at  $\lambda$  differs in  $u'$  and  $v'$ . In  $\underline{M}$  the origin node of  $\lambda$  is  $\mu$ , and in  $\underline{N}$  the origin node of  $\lambda$  is  $\mu'$ . Thus,  $\mu$  and  $\mu'$  are also the origin node of  $\lambda$  in  $M$  and  $N$  respectively for all  $t$  such that  $s \in \underline{t}$ . As we did previously, we first note that the labels at position  $\lambda$  in  $u'$  and  $v'$  cannot be the same constant  $c$ .

If  $\mu \neq \mu'$ , then in  $M$  and  $N$ , for an input tree  $t$  such that  $s \in \underline{t}$ , the origins of  $\lambda$  are different, and the origin functions are not both constants of same image, as the labels at position  $\lambda$  in  $u'$  and  $v'$  cannot be the same constant  $c$ . This means that  $M$  and  $N$  are not weak-origin-equivalent, and thus, by Lemma 19, not equivalent.

If  $\mu = \mu'$ , then the label at  $\lambda$  in  $u'$  and  $v'$  are representations of different functions  $f|_\varphi \neq g|_\psi$ . Since  $M$  and  $N$  are supposed harmonised, we know that  $\varphi = \psi$ . This means that for  $f|_\varphi$  and  $g|_\psi$  to be different, there must exist at least one value  $a \in \varphi$  such that  $f(a) \neq g(a)$ . We pick a tree  $t$  such that  $s \in \underline{t}$  and the label at  $\mu$  in  $t$  is labeled  $a$ . The node at position  $\lambda$  in  $[\underline{M}](t)$  is labeled  $f(a)$ , while the node at position  $\lambda$  in  $[\underline{N}](t)$  is labeled  $g(a)$ . Hence  $[\underline{M}] \neq [\underline{N}]$ .  $\square$

**Corollary 22.** *Under the equivalence-testing assumptions, the equivalence problem for symbolic  $\text{DTOP}$  is decidable.*

*Proof.* Under the equivalence-testing assumptions, the construction of the representation is effective; thus computing the representations and testing their equivalence is an algorithm, as equivalence is decidable for  $\text{DTOP}$ .  $\square$

## 6 Algorithmic Complexity

In this part we study the exact complexity of the algorithm to test equivalence that Corollary 22 points out. Two main points are relevant. First, we note that since our reduction involves a potentially exponential number of various operations related to the equivalence-testing assumptions (computing intersections, deciding function equivalence), their complexity plays an important part in how efficiently symbolic

DTop equivalence can be tested. Second, as previously mentioned, the number of such operations can be made more parsimonious than the method we present in (9)<sub>[p5]</sub>: while the worst case scenario inevitably requires an exponential number of intersections, negations, and emptiness tests on  $\Phi$ , we do not need to build all guards  $\underline{a}$ . Such a construction would only combine all the guards a node can actually encounter in a given symbolic DTop: if a position  $\lambda$  can never be read by a rule of guard  $\varphi$ , we do not care know whether its label satisfies  $\varphi$ , even if  $\varphi$  is otherwise used in the transducer. This construction would be an adaptation of the constuction of the domain automaton of a DTop or the construction of a *compatible* DTop as presented in Lemma 39 and Proposition 62 of [2].

**Corollary 23.** *Under the equivalence-testing assumptions above, the equivalence problem for symbolic DTop is reducible to the equivalence problem on DTop, in  $\text{ExpTime}$  in the worst case, plus, at worst, an exponential number of operations in  $\Phi$  and  $\mathcal{F}$ .*

We note that some realistic cases are more reasonable than this theoretical bound. In the case where guards are all disjoint, the reduction to DTop actually requires a polynomial number of operations in  $\Phi$  and  $\mathcal{F}$ . In practice, it can be expected that few intersections actually need to be computed.

In any case, we can express the number of states and rules in  $\underline{M}$  and  $\underline{N}$  independently of  $\Phi$  and  $\mathcal{F}$ : the states are unchanged, and the number of rules increases, at worst, exponentially:

**Lemma 24.** *For  $M, N$  two symbolic DTop, their DTop representations  $\underline{M}$  and  $\underline{N}$  are DTop with an exponential number of rules and the same number of states.*

Since the problem of DTop equivalence is  $\text{NExpTime}$  [14], a naïve approach to calculating the complexity of symbolic DTop equivalence would yield a  $2\text{NExpTime}$  algorithm, plus an exponential number of operations in  $\Phi$  and  $\mathcal{F}$ . However, upon finer analysis, the complexity of DTop equivalence is tied to the height of a counter-example between two non-equivalent transducers. This height is, in the worst case scenario, exponential in the number of states in the studied DTop. Since representations do not create new states, the height of the counter-examples is unchanged, and the exponentials do not compound.

**Theorem 25.** *The equivalence problem for symbolic DTop is in  $\text{NExpTime}$ , plus, at worst, an exponential number of operations in  $\Phi$  and  $\mathcal{F}$ .*

## 7 Extension to symbolic DTop<sup>R</sup>

We want to extend our results from symbolic DTop to the wider class of symbolic DTop with regular lookahead (symbolic DTop<sup>R</sup>). A symbolic DTop<sup>R</sup> is a pair of a symbolic DTop and its symbolic DTTA lookahead, straight-forwardly extending classical

DTop<sup>R</sup> [9] to the symbolic case. Note that this is different from an inspection; a DTop<sup>R</sup>, whether symbolic or not, has access to the lookahead states in its rules, while a DTop with inspection does not. This means that if  $(M, A)$  is a symbolic DTop<sup>R</sup> describing a transformation of  $\mathcal{T}_\Sigma$ , and  $A$  has a set of states  $Q$ , then  $M$  works on an input signature  $\Sigma \times Q$ .

The class of symbolic DTop<sup>R</sup> is relevant for several reasons. The first is that it is more expressive than the class of symbolic DTop with inspection, and subsumes other relevant classes, such as single-valued symbolic DTop [21, 20]. It also possesses interesting properties. Notably, just as it was the case for DTop<sup>R</sup>, the class of symbolic DTop<sup>R</sup> is closed under composition.

We want to study the equivalence problem of symbolic DTop<sup>R</sup>. For DTop, the addition of a regular lookahead does not prevent the equivalence problem from being decidable, as it is polynomially reduced to equivalence on plain DTop (see [18]). This result can be carried over to the symbolic case, with the same method: annotating the input trees of both symbolic DTop with the states of both lookaheads.

**Lemma 26.** *One can polynomially reduce the equivalence problem of symbolic DTop<sup>R</sup> to the equivalence problem of symbolic DTop with inspection.*

*Proof.* We follow the same method as in the DTop case: if  $N = (M, A)$  and  $N' = (M', A')$  are two symbolic DTop<sup>R</sup>, we want to annotate the trees with the information of both  $A$  and  $A'$ , then test the equivalence of  $M$  and  $M'$  on those trees. The first step is to ensure that the annotation is the combination of the annotation by  $A$  and  $A'$ . This can be done by creating a new lookahead  $A''$  that is the cartesian product of  $A$  and  $A'$ .

As a technical note, the domain of all trees annotated by  $A''$  is regular, but not necessarily recognized by a DTTA. This is a problem, since the equivalence test we produce in this paper for symbolic DTop (as well as the preexisting equivalence test for DTop) requires a domain recognized by a DTTA. This issue is easily fixed by changing the way  $A''$  annotates the input trees: if the information stored at each node is not only the state that  $A''$  reaches but also the rule it uses to do so, then the domain of all annotated trees becomes recognizable by a DTTA.

Finally, testing the equivalence of  $N$  and  $N'$  is exactly testing the equivalence of  $M$  and  $M'$  on all trees annotated by  $A''$ ; note that restricting a symbolic DTop  $M$  to the domain of a DTTA  $A''$  is a simple matter of taking the cartesian product between the states of  $M$  and  $A''$ .  $\square$

We can note that all the steps of this reduction (computing the product of both lookaheads, switching from state to rule annotations, adapting  $M$  and  $M'$  to work on the proper domain) are polynomial, which makes the reduction from symbolic DTop<sup>R</sup> equivalence to symbolic DTop equivalence polynomial. This reduction can be combined with our previous results (Cor. 23 and Thm. 25) to provide the following complexity results:

**Corollary 27.** Under the equivalence-testing assumptions above, the equivalence problem for symbolic  $\text{DTOP}^R$  is reducible to the equivalence problem on  $\text{DTOP}$ , in  $\text{ExpTIME}$  in the worst case, plus, at worst, an exponential number of operations in  $\Phi$  and  $\mathcal{F}$ .

**Theorem 28.** Under the equivalence-testing assumptions, the equivalence problem for symbolic  $\text{DTOP}^R$  is decidable in  $\text{NExpTIME}$ , plus, at worst, an exponential number of operations in  $\Phi$  and  $\mathcal{F}$ .

This result is quite useful, as several  $\text{DTOP}$  classes are fragments of the class of  $\text{DTOP}^R$ . Notably, any  $\text{DBUP}$ , nondeterministic functional  $\text{TOP}$ , or extended  $\text{DTOP}$  can be constructively expressed as  $\text{DTOP}^R$ . Most of those inclusions stand in the symbolic case. The notable exception is extended symbolic  $\text{DTOP}$ : the fact that such transducers can do multi-level pattern matching allows them to manipulate input labels in a way symbolic  $\text{DTOP}^R$  cannot (e.g. switching the labels of parent and child nodes). However, symbolic  $\text{DBUP}$  and nondeterministic symbolic functional  $\text{TOP}$  can be expressed as symbolic  $\text{DTOP}^R$ .

**Corollary 29.** Under the equivalence-testing assumptions, the equivalence problem for deterministic symbolic bottom-up tree transducers and nondeterministic symbolic functional top-down tree transducers is decidable.

## 8 Conclusion

The algorithm presented here provides a novel approach to deciding equivalence for symbolic  $\text{DTOP}$ , and supports non-linear symbolic  $\text{DTOP}$ , by reduction to  $\text{DTOP}$  equivalence. Note that decidability of equivalence for  $\text{DTOP}^R$  [18] works in a comparable way: rather than finding a normal form, the two regular lookaheads are “harmonized” into one, then the problem is reduced to  $\text{DTOP}$  equivalence. The methods presented in this paper also apply to symbolic  $\text{DTOP}^R$  without a critical jump in complexity.

Our method does not involve the computation of a normal form, which is a rather classical technique to decide transducer equivalence [4, 12], with applications to learning. It is interesting to see if normal forms could be defined for symbolic  $\text{DTOP}$ . This looks challenging, however, as it seems more general than finding normal forms for  $\text{DTOP}^R$ , which remains an open problem.

A first possible extension of our model would be to allow the lookahead to have registers, i.e. to memorize some data from the bottom of the tree to annotate the upper part of the tree with it. Under reasonable restrictions, it is likely that we might adapt our methods to reduce the equivalence problem for these objects to the same problem on  $\text{DTOP}^R$ , thus providing a decidability result.

Furthermore, we would like to find out whether this kind of reduction can be applied to more general classes of transducers such as macro tree transducers (with linear size increase), for which equivalence is decidable [11]. If so, then the decidability results

can fairly easily be lifted to symbolic generalisations of the class, at the cost of a few exponential blowups in complexity.

As a final mention, the inversion problem is interesting for symbolic transformations on words and trees, and is relevant to the applications we consider.

## References

1. B. S. Baker. Composition of top-down and bottom-up tree transductions. *Information and Control*, 41(2):186–213, 1979.
2. A. Boiret. *Normalization and Learning of Transducers on Trees and Words*. PhD thesis, Lille University, France, 2016.
3. M. Bojańczyk. Transducers with origin information. In *ICALP 2014*, volume 8573 of *LNCS*, pages 26–37. Springer, 2014.
4. C. Choffrut. Minimizing subsequential transducers: a survey. *Theor. Comput. Sci.*, 292(1):131–143, 2003.
5. H. Comon, M. Dauchet, R. Gilleron, C. Löding, F. Jacquemard, D. Lugiez, S. Tison, and M. Tommasi. Tree automata techniques and applications. Available on: <http://www.grappa.univ-lille3.fr/tata>, 2007. release October, 12th 2007.
6. C. C. Elgot and J. E. Mezei. On relations defined by generalized finite automata. *IBM J. Res. Dev.*, 9(1):47–68, Jan. 1965.
7. J. Engelfriet. Bottom-up and top-down tree transformations - A comparison. *Mathematical Systems Theory*, 9(3):198–231, 1975.
8. J. Engelfriet. Top-down tree transducers with regular look-ahead. *Mathematical systems theory*, 10(1):289–303, 1976.
9. J. Engelfriet. Top-down tree transducers with regular look-ahead. *Mathematical Systems Theory*, 10:289–303, 1977.
10. J. Engelfriet. Some open questions and recent results on tree transducers and tree languages. pages 241–286. Academic Press, 1980.
11. J. Engelfriet and S. Maneth. Macro tree translations of linear size increase are MSO definable. *SIAM J. Comput.*, 32(4):950–1006, 2003.
12. J. Engelfriet, S. Maneth, and H. Seidl. Deciding equivalence of top-down XML transformations in polynomial time. *J. Comput. Syst. Sci.*, 75(5):271–286, 2009.
13. J. Engelfriet and H. Vogler. Macro tree transducers. *J. Comput. Syst. Sci.*, 31(1):71–146, 1985.
14. Z. Ésik. On functional tree transducers. In *FCT*, pages 121–127, 1979.
15. Z. Ésik. Decidability results concerning tree transducers II. *Acta Cybern.*, 6(3):303–314, 1983.
16. Z. Fülöp and H. Vogler. Forward and backward application of symbolic tree transducers. *Acta Inf.*, 51(5):297–325, 2014.
17. A. Lemay, S. Maneth, and J. Niehren. A learning algorithm for top-down XML transformations. In *PODS*, pages 285–296. ACM, 2010.
18. S. Maneth. Equivalence problems for tree transducers: A brief survey. In Z. Ésik and Z. Fülöp, editors, *AFL*, volume 151 of *EPTCS*, pages 74–93, 2014.
19. H. Seidl, S. Maneth, and G. Kemper. Equivalence of deterministic top-down tree-to-string transducers is decidable. In *FOCS*, pages 943–962, 2015.

20. M. Veanes and N. Bjørner. Foundations of finite symbolic tree transducers. *EATCS*, 105:141–173, 2011.
21. M. Veanes and N. Bjørner. Symbolic tree transducers. In *PSI*, volume 7162 of *LNCS*, pages 377–393. Springer, 2011.