

Languages Modulo Normalization

Hitoshi Ohsaki¹ and Hiroyuki Seki²

¹ National Institute of Advanced Industrial Science and Technology
ohsaki@ni.aist.go.jp

² Nara Institute of Science and Technology
seki@is.naist.jp

Abstract. We propose a new class of tree automata, called *tree automata with normalization* (TAN). This framework extends equational tree automata, and improved the results of them: recognized tree languages modulo the *idempotency* $f(x, x) = x$ are closed under complement, besides we do not lose an important decidability result. In the paper, first we investigate the closure properties of this class for Boolean operations and the decidability relative to the equational tree automata. Next we consider the relationship to other automata frameworks, in particular, *hedge automata*, which is a class of unranked tree automata. Hedge automata have been recognized in the XML database community as a theoretical basis for modeling the manipulation of semi-structured data. Through the observation about transformations from XML schema to tree automata, we discuss advantages in the expressiveness and complexity of TAN. As an application of our framework, we show an example of XML schema and constraints that can be dealt with by TAN.

Keywords: tree automata modulo axioms, equational rewriting, Boolean closedness, decidability, regularity, XML schema and constraints.

1 Introduction

Tree automata accept trees as word automata accept words. The class of tree automata hence inherits the benefit from word automata that guarantees important bases for automated reasoning, including closure properties for Boolean operations and the positive decidability results. Propertywise there is no doubt in regular tree automata being a generalization of regular word automata. In fact, several verification tools are designed based on tree automata [1, 11, 12].

However, tree automata accept *irregular* words. For instance, the set of trees $t_{i+1} = f(\mathbf{a}, f(t_i, \mathbf{b}))$ and $t_1 = f(\mathbf{a}, \mathbf{b})$ is regular in the sense of tree automata, while $\mathbf{a}^i \mathbf{b}^i$, which is the leaves of t_i , are no longer regular in words. This is a natural consequence from the fact that trees accepted by regular tree automata are the derivation trees of context-free grammar, and thus, PUMPING LEMMA [4] for regular tree automata can be viewed as a variant of *uvwxy*-lemma for context-free grammar, where for every context-free grammar \mathcal{G} , there exists a natural number k such that, whenever a word z whose length $|z|$ is greater than k is generated by \mathcal{G} , z can be decomposed to be *uvwxy* containing non-empty words v, x and for every $n \geq 0$, $uv^nwx^n y$ is generated by \mathcal{G} .

In the setting of equational tree automata, this becomes much more clear. We showed in [22] that the class of regular tree automata modulo associativity axioms is closely related to the class of context-free grammar. Moreover, *monotone* tree automata modulo associativity which are the super-class of regular tree automata modulo associativity are context-sensitive. These observations are obtained from the property that regular tree automata modulo associativity (for short, regular A-TA) are *not* closed under intersection and complement, and also from that the class of monotone tree automata modulo associativity (monotone A-TA) is *not* decidable in the emptiness, universality and inclusion problems. Furthermore, looking at the transition rules, regular rules $f(\alpha_1(x_1), \dots, \alpha_n(x_n)) \rightarrow \beta(f(x_1, \dots, x_n))$ in tree automata are the production rules $\beta \rightarrow \alpha_1 \cdots \alpha_n$ of context-free grammar, and monotone rules $f(\alpha_1(x_1), \dots, \alpha_n(x_n)) \rightarrow f(\beta_1(x_1), \dots, \beta_n(x_n))$ in tree automata are the rules $\beta_1 \cdots \beta_n \rightarrow \alpha_1 \cdots \alpha_n$ of context-sensitive grammar.

What should be then the counterpart of regular word languages in tree languages? In the paper, we introduce a new tree automata framework, called *tree automata with normalization*, that extends equational tree automata and generalizes the results of them. This class of tree automata includes regular tree automata with associativity normalization (regular TA + R_A) that accept the counterpart of regular word languages.

In equational tree automata, depending on the equational theory, we may lose the Boolean closedness or an important decidability result. In addition to the above-mentioned problems, it is known that regular tree automata modulo ACI theory (*associativity, commutativity and idempotency*) are not closed under complement. There is the same problem in the classes of tree automata modulo *exclusive-or* and *Abelian group* theories. In the left-table below, we summarize the closure properties. The class of languages accepted by *regular tree automata with associativity normalization* is closed under associativity ($=_A$), union (\cup), intersection (\cap) and complement ($(\)^c$). Moreover, the membership ($\in?$), emptiness ($=\emptyset?$), universality ($=T?$) and inclusion ($\subseteq?$) problems are decidable. In the tables, the positive results are indicated by \checkmark , and the negative results are $-$. These results for regular TA + R_A and for other useful classes are the consequence of Theorems 1–5 in the following Sections 3 and 4.

		$=_A$	\cup	\cap	$(\)^c$	$\in?$	$=\emptyset?$	$=T?$	$\subseteq?$
monotone	A-TA (CSG)	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	$-$	$-$	$-$
regular	A-TA (CFG)	\checkmark	\checkmark	$-$	$-$	\checkmark	\checkmark	$-$	$-$
regular	TA + R_A (RG)	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark

The class of the above regular TA + R_A is closely related to *hedge automata*. Hedge automata, originated from [27], are the automata which accepts *unranked* trees, whose transition rules have special patterns which matches regular sequence of state symbols. In practice, every hedge automaton can be minimized,

determinized and completely defined. Moreover, though the instances of transition rules are essentially not finite, most of the decidability questions are solvable. In Section 5, we investigate the relationship between the two frameworks.

As an early work of bi-directional translation between ranked and unranked tree languages, we should address the work of Carme, Niehren and Tommasi [3]. They showed that the hedges over the alphabet Σ are manipulated by a special class of regular tree automata, called *stepwise tree automata*, whose signature is constants from Σ and the binary symbol $@$. In our tree automata framework, their interpretation of hedges and the recognition in stepwise automata are performed by normalization and transition. The detailed discussion on related work and the applications of our framework are summarized in Section 6.

2 Preliminaries

We assume the reader is familiar with term rewriting [2] and tree automata [4]. An *equational theory* is a pair $\mathcal{E} = (F, E)$ in which F is a finite set of function symbols, each with an associated *arity*, and E is a finite set of (orientation sensitive) axioms over the function symbols in F with possibly some variables from V . The binary relation $\rightarrow_{\mathcal{E}}$ induced by \mathcal{E} is the rewrite relation, i.e. $s \rightarrow_{\mathcal{E}} t$ if there is an axiom $l = r$ in E , a context C and a substitution σ such that $s = C[l\sigma]$ and $t = C[r\sigma]$. The equivalence closure of $\rightarrow_{\mathcal{E}}$ is denoted by $=_{\mathcal{E}}$. Associative and commutative theory is an equational theory whose axioms are the associativity and commutativity for some of the binary function symbols. Given a binary function symbol $f \in F$, $f(f(x, y), z) = f(x, f(y, z))$ is an associativity (A) axiom, and $f(x, y) = f(y, x)$ is a commutativity (C) axiom.

An *equational rewriting system* is a pair $\mathcal{R} = (\mathcal{E}, R)$ in which R is a finite set of rewrite rules over F with V . The binary relation $\rightarrow_{\mathcal{R}}$ is a rewrite relation modulo axioms, i.e. $s \rightarrow_{\mathcal{R}} t$ if there is a rewrite rule $l \rightarrow r$ in R , a context C and a substitution σ such that $s =_{\mathcal{E}} C[l\sigma]$ and $t =_{\mathcal{E}} C[r\sigma]$. The reflexive transitive closure and the transitive closure of $\rightarrow_{\mathcal{R}}$ are denoted by $\rightarrow_{\mathcal{R}}^*$ and $\rightarrow_{\mathcal{R}}^+$, respectively. Given a term s , $s \rightarrow_{\mathcal{R}}^! t$ if $s \rightarrow_{\mathcal{R}}^* t$ and t is a normal form. We write $(s) \downarrow_{\mathcal{R}}^!$ to denote such t . For the set $\mathcal{T}(F)$ of ground terms over the signature F , $\text{NF}_F(\mathcal{R})$ represents the set of normal forms of \mathcal{R} over F . An equational rewrite system \mathcal{R} is weakly normalizing if for every term t , there exists some $(t) \downarrow_{\mathcal{R}}^!$. We say \mathcal{R} (R modulo \mathcal{E}) is confluent/strongly normalizing if $\rightarrow_{\mathcal{R}}$ is confluent/strongly normalizing.

Given an equational theory $\mathcal{E} = (F, E)$ and an equational rewriting system $\mathcal{R} = (\mathcal{E}, R)$, the \mathcal{E} -closure of a subset L of terms, denoted by $\mathcal{E}(L)$, is the smallest set containing L closed under $=_{\mathcal{E}}$, i.e. if $s =_{\mathcal{E}} t$ and $s \in \mathcal{E}(L)$ then $t \in \mathcal{E}(L)$. The \mathcal{R} -descendant closure of L , $\mathcal{R}^*(L)$, is the smallest set containing L closed under $\rightarrow_{\mathcal{R}}^*$. Similarly, the \mathcal{R} -ascendant closure $(\mathcal{R}^{-1})^*(L)$ is the smallest set containing L closed under $\rightarrow_{\mathcal{R}^{-1}}^*$ and the \mathcal{R} -equivalence closure $(\mathcal{R} \cup \mathcal{R}^{-1})^*(L)$ is the smallest set containing L closed under $=_{\mathcal{R}}$.

A *tree automaton with normalization* (TAN) is a tuple $\mathcal{A} = (\mathcal{R}, Q, Q_{\text{fin}}, \Delta)$, consisting of the equational rewrite system $\mathcal{R} = (\mathcal{E}, R)$ with $\mathcal{E} = (F, E)$, a finite

set Q of states disjoint from the symbols in F , a subset Q_{fin} of Q , and a finite set Δ of transition rules whose shapes are in the following forms:

$$\begin{array}{ll} \text{(REGULAR)} & \text{(EPSILON)} \\ f(\alpha_1(x_1), \dots, \alpha_n(x_n)) \rightarrow \beta(f(x_1, \dots, x_n)) & \alpha(x_1) \rightarrow \beta(x_1) \end{array}$$

for some $f \in F$ with $\text{arity}(f) = n$ and $\alpha_1, \dots, \alpha_n, \alpha, \beta \in Q$. Every state is a unary symbol, and variables x_i, x_j in the transition rule are different if $i \neq j$.

A *move relation* of \mathcal{A} is concatenation of the two kinds of relations over $\mathcal{T}(F \cup Q)$, the normalization $\rightarrow_{\mathcal{R}}^!$ and the transition \rightarrow_{Δ} modulo \mathcal{E} , that means: $s \rightarrow_{\mathcal{A}} t$ if (1) there is a normal form u of s with respect to $\rightarrow_{\mathcal{R}}$ and (2) there is a transition rule $l \rightarrow r$ in Δ , a context C and a substitution σ such that $u =_{\mathcal{E}} C[l\sigma]$ and $t =_{\mathcal{E}} C[r\sigma]$. If there is no normal form $(s) \downarrow_{\mathcal{R}}^!$, the transition fails.

A tree t is *accepted* by \mathcal{A} if $t \in \mathcal{T}(F)$ and $t \rightarrow_{\mathcal{A}}^* \alpha(t')$ for some $\alpha \in Q_{\text{fin}}$ and $t' \in \mathcal{T}(F)$. The set of trees accepted by \mathcal{A} is denoted by $\mathcal{L}(\mathcal{A})$. We say a TAN $\mathcal{A} = (\mathcal{R}, Q, Q_{\text{fin}}, \Delta)$ with $\mathcal{R} = (F, E, R)$ is *regular* if transition rules in Δ are all in the REGULAR form. In case of $R = \emptyset$ and $E = \emptyset$, regular TAN are called regular tree automata. A tree language accepted by a regular tree automaton is called *regular*. These notions of regularity coincide with, if $R = \emptyset$ and $E = \emptyset$, the standard definition, e.g. [4].

As an example, we consider the TAN \mathcal{A}_1 with

$$\begin{array}{l} R_1: f(\mathbf{h}(x), \mathbf{h}(x)) \rightarrow \mathbf{h}(x) \\ E_1: f(f(x, y), z) = f(x, f(y, z)) \quad f(x, y) = f(y, x) \\ \Delta_1: \mathbf{a} \rightarrow \alpha(\mathbf{a}) \quad \mathbf{g}(\alpha(x)) \rightarrow \alpha(\mathbf{g}(x)) \quad \mathbf{h}(\alpha(x)) \rightarrow \beta(\mathbf{h}(x)) \end{array}$$

where β is the final state. Then $\mathbf{h}(\mathbf{g}(\mathbf{a}))$ is accepted by \mathcal{A}_1 , because $\mathbf{h}(\mathbf{g}(\mathbf{a})) \rightarrow_{\mathcal{R}_1}^! \mathbf{h}(\mathbf{g}(\mathbf{a}))$ and $\mathbf{h}(\mathbf{g}(\mathbf{a})) \rightarrow_{\Delta_1} \mathbf{h}(\mathbf{g}(\alpha(\mathbf{a})))$, and this turns out $\mathbf{h}(\mathbf{g}(\mathbf{a})) \rightarrow_{\mathcal{A}_1}^* \beta(\mathbf{h}(\mathbf{g}(\mathbf{a})))$.

Regarding the language accepted by this example, we observe that

if a tree t consists of nodes f with the same tree $\mathbf{h}(\mathbf{g}^i(\mathbf{a}))$ at each leaf ($i \geq 0$),
 t is accepted by \mathcal{A}_1 ,

because for every such tree t , $t \rightarrow_{\mathcal{R}_1}^! \mathbf{h}(\mathbf{g}^i(\mathbf{a}))$, and then $\mathbf{h}(\mathbf{g}^i(\mathbf{a})) \rightarrow_{\mathcal{A}_1}^* \beta(\mathbf{h}(\mathbf{g}^i(\mathbf{a})))$. For instance, $f(\mathbf{h}(\mathbf{g}(\mathbf{a})), \mathbf{h}(\mathbf{g}(\mathbf{a}))) \rightarrow_{\mathcal{R}_1}^! \mathbf{h}(\mathbf{g}(\mathbf{a}))$, and thus $\mathbf{h}(\mathbf{g}(\mathbf{a})) \rightarrow_{\mathcal{A}_1}^* \beta(\mathbf{h}(\mathbf{g}(\mathbf{a})))$ from the above example.

In the next example, we compare the expressiveness between TAN and ETA. ETA (equational tree automata [23]) are the special class of TAN whose set R of rewrite rules is empty. We take the TAN \mathcal{A}_2 with the following components:

$$\begin{array}{l} R_2: f(\mathbf{h}(x), \mathbf{h}(x)) \rightarrow \mathbf{h}(x) \\ E_2: f(f(x, y), z) = f(x, f(y, z)) \quad f(x, y) = f(y, x) \\ \Delta_2: \mathbf{a} \rightarrow \alpha(\mathbf{a}) \quad \mathbf{g}(\alpha(x)) \rightarrow \alpha(\mathbf{g}(x)) \quad \mathbf{h}(\alpha(x)) \rightarrow \beta(\mathbf{h}(x)) \\ \quad f(\beta(x), \beta(y)) \rightarrow \gamma(f(x, y)) \end{array}$$

In this example, \mathcal{A}_2 has the the additional transition rule $f(\beta(x), \beta(y)) \rightarrow \gamma(f(x, y))$ with the final state γ . Similar to the previous example, we observe that

if a tree t consists of nodes f with two different kinds of trees $h(g^i(a))$ and $h(g^j(a))$ at its leaves ($i \neq j$), t is accepted by \mathcal{A}_2 .

For instance, $f(h(a), f(h(g(a)), h(a)))$ is accepted: $f(h(a), f(h(g(a)), h(a))) \xrightarrow{\dagger_{\mathcal{R}_2}} f(h(a), h(g(a)))$, and therefore $f(h(a), f(h(g(a)), h(a))) \xrightarrow{*_{\mathcal{A}_2}} \gamma(f(h(a), h(g(a))))$. However, $f(h(a), h(a))$ is not accepted by \mathcal{A}_2 , because $f(h(a), h(a)) \xrightarrow{\dagger_{\mathcal{R}_2}} h(a)$ and $h(a) \xrightarrow{*_{\mathcal{A}_2}} \beta(h(a))$. In fact, $f(h(a), h(a)) \xrightarrow{*_{\mathcal{A}_2}} \beta(h(a))$.

Remark 1. The previous language $\mathcal{L}(\mathcal{A}_2)$ can not be represented by ETA if $E = \{f(h(x), h(x)) = h(x), f(f(x, y), z) = f(x, f(y, z)), f(x, y) = f(y, x)\}$.

Furthermore, the language $\mathcal{L}(\mathcal{A}_2)$ is not accepted by multitree automata [14] or two-way equational tree automata [28] either.

3 Relative Decidability

We begin with the membership problem for the class of TAN. Hereafter in the following two sections, we suppose $\mathcal{E} = (F, E)$ and $\mathcal{R} = (\mathcal{E}, R)$.

Theorem 1. *The membership problem for \mathcal{A} is decidable if \mathcal{R} is strongly normalizing over $\mathcal{T}(F)$ and the membership problem for an equational tree automaton \mathcal{A}' is decidable, where \mathcal{A}' is defined by replacing R in \mathcal{A} with \emptyset . \square*

In this framework, the decidability results depend upon that of equational tree automata, because we have $(\xrightarrow{\dagger_{\mathcal{R}}} \cdot \xrightarrow{\dagger_{\mathcal{A}'}})^+ \subseteq \xrightarrow{\dagger_{\mathcal{R}}} \cdot \xrightarrow{\dagger_{\mathcal{A}'}}$. To be precise, we state the following lemma.

Lemma 1. *Given $\mathcal{A} = (\mathcal{R}, Q, Q_{\text{fin}}, \Delta)$, let $\mathcal{A}' = (\mathcal{R}', Q, Q_{\text{fin}}, \Delta)$ with $\mathcal{R}' = (\mathcal{E}, \emptyset)$, then for every $t \in T(F)$, $t \xrightarrow{\dagger_{\mathcal{A}}} t'$ implies $(t) \downarrow_{\mathcal{R}}^{\dagger} \xrightarrow{*_{\mathcal{A}'}} t'$. \square*

Accordingly, the question if $t \in \mathcal{L}(\mathcal{A})$ is determined by testing $(t) \downarrow_{\mathcal{R}}^{\dagger} \in \mathcal{L}(\mathcal{A}')$ for all $(t) \downarrow_{\mathcal{R}}^{\dagger}$. From König's Lemma, when \mathcal{R} is strongly normalizing, there are only finitely many candidates for $(t) \downarrow_{\mathcal{R}}^{\dagger}$. Strong normalization in Theorem 1, however, can not be weakened to weak normalization (Appendix A).

Every language accepted by a TAN $\mathcal{A} = (\mathcal{R}, Q, Q_{\text{fin}}, \Delta)$ is the \mathcal{R} -ascendant closure of some subset of normal forms: Let $\mathcal{A}' = ((\mathcal{E}, \emptyset), Q, Q_{\text{fin}}, \Delta)$,

$$(i) \quad \mathcal{L}(\mathcal{A}) = (\mathcal{R}^{-1})^*(\text{NF}_F(\mathcal{R}) \cap \mathcal{L}(\mathcal{A}')).$$

Moreover,

$$(ii) \quad \mathcal{E}(\mathcal{L}(\mathcal{A})) = \mathcal{L}(\mathcal{A}),$$

$$(iii) \quad (\mathcal{R} \cup \mathcal{R}^{-1})^*(\mathcal{L}(\mathcal{A})) = \mathcal{R}^*(\mathcal{L}(\mathcal{A})) = \mathcal{L}(\mathcal{A}) \text{ if } \mathcal{R} \text{ is confluent.}$$

For instance, we consider the TAN \mathcal{A}_1 with

$$R: f(h(x), h(x)) \rightarrow h(x)$$

$$E: f(f(x, y), z) = f(x, f(y, z)) \quad f(x, y) = f(y, x)$$

$$\Delta: a \rightarrow \alpha(a) \quad g(\alpha(x)) \rightarrow \alpha(g(x)) \quad h(\alpha(x)) \rightarrow \beta(h(x))$$

We have $\mathcal{L}(\mathcal{A}_1) = (\mathcal{R}^{-1})^*(\mathcal{L}(\mathcal{A}'_1))$. This follows from Property (i) and $\text{NF}_F(\mathcal{R}) \cap \mathcal{L}(\mathcal{A}'_1) = \mathcal{L}(\mathcal{A}'_1)$. Furthermore, $(\mathcal{R} \cup \mathcal{R}^{-1})^*(\mathcal{L}(\mathcal{A}'_1)) = \mathcal{L}(\mathcal{A}_1)$, because $\mathcal{L}(\mathcal{A}'_1) \subseteq \mathcal{L}(\mathcal{A}_1)$ from $\mathcal{L}(\mathcal{A}_1) = (\mathcal{R}^{-1})^*(\mathcal{L}(\mathcal{A}'_1))$. Since \mathcal{R} is confluent, $(\mathcal{R} \cup \mathcal{R}^{-1})^*(\mathcal{L}(\mathcal{A}'_1)) = \mathcal{L}(\mathcal{A}_1)$ from Property (iii). Therefore, $(\mathcal{R}^{-1})^*(\mathcal{L}(\mathcal{A}'_1)) \subseteq (\mathcal{R} \cup \mathcal{R}^{-1})^*(\mathcal{L}(\mathcal{A}'_1)) \subseteq (\mathcal{R} \cup \mathcal{R}^{-1})^*(\mathcal{L}(\mathcal{A}_1))$.

Theorem 2. *The emptiness problem is decidable for tree automata with normalization if*

- (1) *there effectively exists an equational tree automaton (with the same equational theory) which accepts the set of normal forms of \mathcal{R} ,*
- (2) *this class of equational tree automata is closed under the intersection.*

Proof. Let $\mathcal{A}_R = ((\mathcal{E}, R), P, P_{\text{fin}}, \Delta)$ and $\mathcal{A}_\emptyset = ((\mathcal{E}, \emptyset), P, P_{\text{fin}}, \Delta)$. The ETA \mathcal{A}_\emptyset is obtained from the TAN \mathcal{A}_R by replacing R with \emptyset . From the assumption (1), we obtain an ETA \mathcal{B} , equipped with the same theory $\mathcal{E} = (F, E)$, that accepts $\text{NF}_F(\mathcal{R})$. This implies from Property (i) that

$$\mathcal{L}(\mathcal{A}_R) = \emptyset \Leftrightarrow \mathcal{L}(\mathcal{A}_\emptyset) \cap \mathcal{L}(\mathcal{B}) = \emptyset.$$

From the assumption (2), we have an ETA $\mathcal{C} = ((F, E, \emptyset), Q, Q_{\text{fin}}, \Lambda)$ that accepts $\mathcal{L}(\mathcal{A}_\emptyset) \cap \mathcal{L}(\mathcal{B})$. From Lemma 9 (Appendix A), $\mathcal{L}(\mathcal{C})$ is the \mathcal{E} -closure of the language accepted by $\mathcal{C}' = ((F, \emptyset, \emptyset), Q, Q_{\text{fin}}, \Lambda)$. Hence, $\mathcal{L}(\mathcal{A}_R) = \emptyset \Leftrightarrow \mathcal{L}(\mathcal{C}') = \emptyset$. Because \mathcal{C}' is in the class of regular tree automata, the emptiness question to \mathcal{C}' is decidable. \square

The first equivalence in the above proof can be rephrased as follows: $\mathcal{L}(\mathcal{A}_R) \neq \emptyset$ if and only if $\mathcal{L}(\mathcal{A}_\emptyset)$ contains a normal form of \mathcal{R} . When $E = \emptyset$ in Theorem 2, the latter question is known to be EXPTIME-complete regardless of R [5].

Corollary 1. *If $E = \emptyset$, the conditions (1) and (2) can be omitted.* \square

Theorem 3. *The inclusion problem is decidable for tree automata with normalization if the conditions (1),(2) hold and the inclusion problem is decidable for this class of equational tree automata.* \square

It is known that, if \mathcal{E} is AC theory for a binary symbol f , the inclusion problem for ETA is decidable. Even if \mathcal{E} also contains the axioms of $f(x, 0) = x$ and/or $g(f(x, y)) = f(g(x), g(y))$, the inclusion problem is decidable [28, 29].

Let \mathcal{B} and \mathcal{D} be ETA equipped with the equational theory \mathcal{E} , such that \mathcal{B} accepts $\text{NF}_F(\mathcal{R})$ and \mathcal{D} accepts trees which are not normalizable. The universality question if $\mathcal{L}(\mathcal{A}_R) = \mathcal{T}(F)$ is equivalent to ask if $\mathcal{L}(\mathcal{B}) \subseteq \mathcal{L}(\mathcal{A}_\emptyset)$ and $\mathcal{L}(\mathcal{D}) = \emptyset$:

Theorem 4. *The universality is decidable for tree automata with normalization if the conditions (1),(3) hold and the inclusion problem is decidable for this class of equational tree automata:*

- (3) *there effectively exists an equational tree automaton (with the same equational theory) which accepts the set of non-normalizable trees.* \square

Corollary 2. *If \mathcal{R} is weakly normalizing, the condition (3) can be omitted.* \square

4 Closedness under Boolean Operations

Lemma 2. *For every equational rewrite system \mathcal{R} , the class of tree automata with \mathcal{R} -normalization is closed under complement if \mathcal{R} is weakly normalizing and confluent over $\mathcal{T}(F)$ and this class of equational tree automata is closed under complement. \square*

In the previous lemma, weak normalization is essential. Let us consider $\mathcal{R}_1 = ((\{a, b\}, \emptyset), \{a \rightarrow a\})$. Because a is not normalized, by definition, a is not accepted by TAN with \mathcal{R}_1 . This implies that though the language $\{b\}$ is accepted by $\mathcal{A}_1 = (\mathcal{R}_1, \{\alpha\}, \{\alpha\}, \{b \rightarrow \alpha(b)\})$, the complement of $\mathcal{L}(\mathcal{A}_1)$ is not accepted by TAN with \mathcal{R}_1 .

Also, confluence is necessary for closure under complement. We consider $\mathcal{R}_2 = ((\{a, b, c\}, \emptyset), \{a \rightarrow b, a \rightarrow c\})$ and $\mathcal{A}_2 = (\mathcal{R}_2, \{\alpha\}, \{\alpha\}, \{b \rightarrow \alpha(b)\})$. Because a is normalized to b , a is accepted by \mathcal{A}_2 as b is accepted. Every TAN with \mathcal{R}_2 accepting c accepts a , because $a \rightarrow_{\mathcal{R}_2}^! c$. This implies that the complement of $\mathcal{L}(\mathcal{A}_2)$ is not accepted by TAN with \mathcal{R}_2 .

Theorem 5. *For every equational rewrite system \mathcal{R} , the class of tree automata with \mathcal{R} -normalization is closed under Boolean operations if \mathcal{R} is weakly normalizing and confluent over $\mathcal{T}(F)$ and this class of equational tree automata is closed under Boolean operations. \square*

In the above theorem, confluence is essential for closure under intersection: Consider the previous example $\mathcal{R}_2 = ((\{a, b, c\}, \emptyset), \{a \rightarrow b, a \rightarrow c\})$. Let $\mathcal{A}_2 = (\mathcal{R}_2, \{\alpha\}, \{\alpha\}, \{b \rightarrow \alpha(b)\})$ and $\mathcal{A}_3 = (\mathcal{R}_2, \{\alpha\}, \{\alpha\}, \{c \rightarrow \alpha(c)\})$. The intersection of the two languages $\mathcal{L}(\mathcal{A}_2)$ and $\mathcal{L}(\mathcal{A}_3)$, which is $\{a, b\} \cap \{a, c\}$, is not accepted by TAN with \mathcal{R}_2 , because every TAN with \mathcal{R}_2 accepting a accepts b or c .

However, weak normalization is not needed for closure under intersection. Moreover, weak normalization and confluence are not required for the union.

Lemma 3. *Let $\mathcal{E} = (F, AC)$ whose component AC is the set of AC axioms for f . If $R \in \{R_U, R_I, R_N, R_D, R_{UI}, R_{UN}, R_{UD}, R_{AG}\}$ defined below, (\mathcal{E}, R) is strongly normalizing and confluent over $\mathcal{T}(F)$:*

$$R_U: f(x, 0) \rightarrow x \quad R_{UI}: R_U \cup R_I \quad R_{UN}: R_U \cup R_N \quad R_{UD}: R_U \cup R_D$$

$$R_I: f(x, x) \rightarrow x$$

$$R_N: f(x, x) \rightarrow 0$$

$$R_D: g(f(x, y)) \rightarrow f(g(x), g(y)) \quad g(0) \rightarrow 0 \quad g(g(x)) \rightarrow x$$

$$R_G: R_{UD} \cup \{f(x, g(x)) \rightarrow 0\} \quad (\text{ABELIAN GROUP})$$

Proof. Strong normalization can easily be shown, and thus, confluence follows from CRITICAL PAIR LEMMA [2]. Tool support, e.g. that of CiME [6], is helpful in checking all AC critical pairs to be joinable. \square

Corollary 3. *If $R \in \{R_U, R_I, R_N, R_D, R_{UI}, R_{UN}, R_{UD}, R_G\}$ and E is the set of AC axioms for f , the class of tree automata with normalization is closed under Boolean operations. Moreover, the membership problem is decidable for this class of tree automata with normalization. \square*

Remark 2. Equational tree automata with ACUI (resp. ACUN) axioms are *not* closed under complement. Applying a similar proof of Remark 1 (Appendix A), we can even show that equational tree automata with ACI (resp. ACN) axioms are also *not* closed under complement. The same problem is discussed in [29].

Furthermore, it is not obvious for ETA how to deal with Boolean operations in ACD and ACUD cases. Though it is shown in [29] that the class of tree languages modulo ACUD is closed under Boolean operations, the computation algorithm is quite involved. In our framework, Boolean computation for TAN with R_{UD} modulo AC can be manipulated by the operations for AC tree automata.

Lemma 4. *There effectively exists an AC-tree automaton which accepts the set of normal forms of R_U, R_D, R_{UD} modulo the AC axioms, respectively.*

Proof. We consider R_U . Because of left-linearity, there effectively exists a tree automaton $\mathcal{A} = ((F, \emptyset, \emptyset), Q, Q_{\text{fin}}, \Delta)$ which accepts normal forms of R_U . The set of normal forms of R_U modulo AC is the AC-closure of $\mathcal{L}(\mathcal{A})$, which is accepted by $((F, \text{AC}, \emptyset), Q, Q_{\text{fin}}, \Delta)$. Similarly, we can prove for R_D and R_{UD} . \square

Corollary 4. *The emptiness, universality and inclusion problems are decidable for R_U, R_D, R_{UD} modulo the AC axioms. \square*

Corollary 5. *The emptiness problem is decidable for R_I, R_N, R_{UI}, R_{UN} modulo the AC axioms if the following question is decidable: Given a regular tree automaton \mathcal{A} , does its language contain a normal form of R_I modulo AC? \square*

5 Regular Leaf Languages and Hedge Automata

When discussing the expressiveness of tree languages, we often take the images of them into word languages. Leaves of a tree t are the sequence in the left-to-right order of constants occurring in t . We denote leaf for the mapping that takes a tree and returns the leaves, i.e. $\text{leaf}(f(t_1, \dots, t_n)) = \text{leaf}(t_1) \cdots \text{leaf}(t_n)$ if $n \geq 1$; $\text{leaf}(a) = a$, otherwise. The leaf language associated to a tree language L is the set of leaves obtained from L .

Given a signature F , we take the rewrite system \mathcal{R}_{reg} whose rewrite rules are for arbitrary function symbols f, g in F ,

$$f(x_1, \dots, g(y_1, y_2, \dots, y_k), x_i, \dots, x_n) \rightarrow f(x_1, \dots, y_1, g(y_2, \dots, y_k, x_i), \dots, x_n)$$

such that $\text{arity}(f) \geq 2$ and $\text{arity}(g) \geq 1$. Here f and g are possibly the same. Obviously, \mathcal{R}_{reg} is strongly normalizing if $E = \emptyset$. However, \mathcal{R}_{reg} is *not* confluent, because $f(g(h(a, b), c), d) \rightarrow_{\mathcal{R}_{\text{reg}}}^! f(a, h(b, g(c, d)))$ and $f(g(h(a, b), c), d) \rightarrow_{\mathcal{R}_{\text{reg}}}^! f(a, g(b, h(c, d)))$, though it preserves the leaves.

Lemma 5. *Suppose that for every s, t in $\mathcal{T}(F)$, $s =_{\mathcal{E}} t$ implies $\text{leaf}(s) = \text{leaf}(t)$. Then, $t =_{\mathcal{R}_{\text{reg}}} u$ if and only if $\text{leaf}(t) = \text{leaf}(u)$. \square*

Trees in \mathcal{R}_{reg} -normal form, if $E = \emptyset$, are *right* associative as derivation trees of regular word grammar are right associative. For instance, when the grammar has the production rules $\alpha \rightarrow \mathbf{a}\beta$, $\beta \rightarrow \mathbf{b}\gamma$, $\gamma \rightarrow \mathbf{c}$, we have $\alpha(\mathbf{a} \circ \beta(\mathbf{b} \circ \gamma(\mathbf{c})))$ for the derivation starting from α . The word \mathbf{abc} is the leaves of this derivation tree. Meanwhile, in tree automata setting, the word \mathbf{abc} can be represented to be $\mathbf{a} \cdot (\mathbf{b} \cdot \mathbf{c})$. Tree automaton accepting $\mathbf{a} \cdot (\mathbf{b} \cdot \mathbf{c})$ is equipped with the transition rules $\mathbf{a} \rightarrow p_{\mathbf{a}}(\mathbf{a})$, $p_{\mathbf{a}}(x) \cdot \beta(y) \rightarrow \alpha(x \cdot y)$, $\mathbf{b} \rightarrow p_{\mathbf{b}}(\mathbf{b})$, $p_{\mathbf{b}}(x) \cdot \gamma(y) \rightarrow \beta(x \cdot y)$, $\mathbf{c} \rightarrow \gamma(\mathbf{c})$. Here we use the infix operators \circ and \cdot for denoting the branching node of derivations and the concatenation of alphabet.

Accordingly, trees which have such *regular* leaves should contain constants that appear on the left in a tree whenever non-constant trees appear at the same level in the tree. If \mathcal{A} is a tree automaton, the leaf language $\{\text{leaf}(t) \mid (t) \downarrow_{\mathcal{R}_{\text{reg}}}^!\}$ accepted by \mathcal{A} is regular, while $\{\text{leaf}(t) \mid t \text{ accepted by } \mathcal{A}\}$ is context-free. Since trees accepted by tree automata are derivation trees generated by context-free grammar, this observation can be easily obtained. But what about a tree automaton accepting $\{t \mid \text{leaf}(t) \text{ generated by } \mathcal{G}\}$?

Lemma 6. *Given a regular grammar $\mathcal{G} = (\Sigma, \mathcal{S}, s_0, A)$, we can construct a regular tree automaton \mathcal{A} which accepts all trees whose leaves are generated by a regular grammar \mathcal{G} . \square*

Corollary 6. *The above lemma can be extended to an arbitrary signature F if F contains constants from Σ and a binary function symbol. \square*

Of particular interest about Lemma 6 is that this property does not hold for context-free grammar. The proof of this observation is found in Appendix A.

Lemma 7. *For every regular tree automaton \mathcal{A} , $\{t \mid (t) \downarrow_{\mathcal{R}_{\text{reg}}}^!\}$ accepted by \mathcal{A} is a regular tree language. \square*

Hence, the class of tree automata, each of which accepts all trees whose leaves is generated by a regular grammar, is a proper sub-class of regular tree automata. In the framework of TAN, this statement can be represented as follows.

Theorem 6. *The class of TAN with \mathcal{R}_{reg} is a sub-class of regular tree automata. Moreover, for every confluent sub-system \mathcal{R} of \mathcal{R}_{reg} , the class of TAN with \mathcal{R} is also a sub-class of regular tree automata, and is*

- closed under $=_{\mathcal{R}}$,
- closed under Boolean operations,
- decidable in the membership, emptiness, universality and inclusion problems. \square

An associativity rewrite system \mathcal{R}_A , whose rewrite rules are $f(f(x, y), z) \rightarrow f(x, f(y, z))$ for some of the binary function symbols, is such an example of confluent sub-systems. This class of rewrite systems is useful in translating *unranked* trees in ranked signature.

In case of *unranked* signature, trees with regular leaves are manipulated by *hedge automata*. Recently, the hedge automata theory is accepted in XML community [20], and in which XML documents are considered as tree-structured objects, called *hedges*. Hence, as if an XML schema matches documents, a hedge automaton accepts the set of hedges. Hedges in $\text{HDG}(\Sigma, C)$ over a finite set Σ of label symbols and a finite set C of constants are recursively defined as follows:

$$\begin{aligned}
t \in \text{HDG}(\Sigma, C) & \quad \text{if } t = c & \quad \text{for some } c \in C \\
& \quad \text{if } t = a\langle w \rangle & \quad \text{for some } a \in \Sigma, w \in \text{HDGS}(\Sigma, C) \\
w \in \text{HDGS}(\Sigma, C) & \quad \text{if } w = \epsilon & \quad (\text{null hedge}) \\
& \quad \text{if } w = t & \quad \text{for some } t \in \text{HDG}(\Sigma, C) \\
& \quad \text{if } w = u \circ v & \quad \text{for some } u, v \in \text{HDGS}(\Sigma, C)
\end{aligned}$$

The infix operator \circ over hedges is associative over hedges with the identity axiom for ϵ . In the unranked setting, $t_1 \circ \dots \circ t_n$ can be denoted as $\circ(t_1, \dots, t_n)$, and it satisfies that $\circ(t_1, \dots, t_i, \epsilon, t_{i+1}, \dots, t_n) = \circ(t_1, \dots, t_i, t_{i+1}, \dots, t_n)$.

One should remark on hedges defined above. In the literature, e.g. [19], the null hedge ϵ and the sequence of hedges are defined as well-formed hedges. However, this leads us to confusion in distinguishing *elements* from *contents* in the context of XML. According to XML schema recommended by W3C [30], an *empty-element tag*, that corresponds to some c in C , is an element but is not the empty sequence of elements. On the other hand, an element (tagged a) with no content means the element having the empty sequence in between *start-tag* and *end-tag*, that corresponds $a\langle \ \rangle$.

A hedge automaton is a tuple $\mathcal{H} = (\Sigma, C, N, N_{\text{fin}}, \Delta)$, in which N is a finite set of unary state symbols disjoint from Σ and C , N_{fin} is a set of final states with $N_{\text{fin}} \subseteq N$, and Δ is a set of transition rules for Σ and C . Transition rules in Δ are in the two shapes of

$$\begin{aligned}
& c \rightarrow \beta(c) \\
& a\langle \alpha_1(x_1) \circ \dots \circ \alpha_n(x_n) \rangle \rightarrow \beta(a\langle x_1 \circ \dots \circ x_n \rangle) \Leftarrow \alpha_1 \dots \alpha_n \in \mathcal{L}(\mathcal{G})
\end{aligned}$$

with $a \in \Sigma$, $x_1, \dots, x_n \in V$ and $\alpha_1, \dots, \alpha_n, \beta \in N$.

Formally, we should introduce *sequence variables* ([13]) for the second type of rules, so the left-hand side can match an arbitrary sequence $\alpha_1(x_1) \circ \dots \circ \alpha_n(x_n)$. In the sense of tree automata, we thus allow hedge automata to contain in Δ infinite number of regular transition rules in a certain type.

The grammar \mathcal{G} over N that appears in the second type of rules must be regular, but that can be different for each transition rule. The left-arrow (\Leftarrow) together with the membership condition represents the guard, meaning that, if the sequence $\alpha_1 \dots \alpha_n$ of states is generated by \mathcal{G} , the transition fires. Variables x_1, \dots, x_n in a rule are instantiated to certain hedges when the rule is applied.

In tree automata framework, hedges can be modeled over the signature

$$\begin{aligned}
F_{\Sigma, C}: \{ a \mid a \in \Sigma \} & \quad \text{arity}(a) = 1 \quad \text{for each } a, & \{ \circ \} & \quad \text{arity}(\circ) = 2, \\
& \{ c \mid c \in C \} & \quad \text{arity}(c) = 0 \quad \text{for each } c, & \{ \epsilon \} & \quad \text{arity}(\epsilon) = 0.
\end{aligned}$$

$$\begin{aligned}
\mathcal{H}_{\text{in}}: \quad & f\langle \gamma(x) \circ \beta(y) \rangle \rightarrow \gamma(f\langle x \circ y \rangle) \\
& g\langle \alpha(z_1) \circ \cdots \circ \alpha(z_m) \rangle \rightarrow \gamma(g\langle z_1 \circ \cdots \circ z_m \rangle) \\
& g\langle \beta(w_1) \circ \cdots \circ \beta(w_n) \rangle \rightarrow \delta(g\langle w_1 \circ \cdots \circ w_n \rangle) \\
& \mathbf{a} \rightarrow \alpha(\mathbf{a}) \\
& \mathbf{b} \rightarrow \beta(\mathbf{b})
\end{aligned}$$

Fig. 1. Example of XML Schema

Hereafter, to discuss the translation from unranked to ranked signatures, we define the *unflattening* operation:

$$\text{unflat}(t) = \begin{cases} a(\text{unflat}(t_1) \circ \cdots (\text{unflat}(t_n) \circ \text{unflat}(t_{n+1})) \cdots) & \text{if } t = a\langle t_1 \cdots t_n t_{n+1} \rangle \text{ with } n \geq 0, \\ a(\epsilon) & \text{if } t = a\langle \rangle, \\ t & \text{if } t \in C. \end{cases}$$

Moreover, we use *flat* for the reverse translation, e.g. $\text{flat}(f((\mathbf{a} \circ \mathbf{b}) \circ \mathbf{g}(\mathbf{c}))) = f\langle \mathbf{a} \circ \mathbf{b} \circ \mathbf{g}(\mathbf{c}) \rangle$.

Proposition 1. *For every hedge automaton \mathcal{H} over Σ and C , we can construct a tree automaton $\mathcal{A}_{\mathcal{H}}$ over $F_{\Sigma, C}$ such that t is accepted by \mathcal{H} if and only if $\text{unflat}(t)$ is accepted by $\mathcal{A}_{\mathcal{H}}$. \square*

In the above setting, $f(\mathbf{a} \circ (\mathbf{b} \circ \mathbf{c}))$ and $f((\mathbf{a} \circ \mathbf{b}) \circ \mathbf{c})$ are distinguished, because we *do not* assume that $(x \circ y) \circ z = x \circ (y \circ z)$ is given. However, due to Lemma 7, we can generalize Proposition 1 as follows.

Corollary 7. *For every hedge automaton \mathcal{H} over Σ and C , we can construct a tree automaton $\mathcal{B}_{\mathcal{H}}$ over $F_{\Sigma, C}$ such that $\text{flat}(t)$ is accepted by \mathcal{H} if and only if t is accepted by $\mathcal{B}_{\mathcal{H}}$. \square*

The sizes $|\mathcal{A}_{\mathcal{H}}|$ and $|\mathcal{B}_{\mathcal{H}}|$ are $O(|\mathcal{H}|)$ and $O(|\mathcal{H}|^3)$, respectively. The time complexities are linear for the construction of $\mathcal{A}_{\mathcal{H}}$, and cubic for $\mathcal{B}_{\mathcal{H}}$ to the size of \mathcal{H} . In Corollary 7, we can take the TAN $\mathcal{A}_{\mathcal{H}}$ with $\mathcal{R}_{\mathbf{A}}$, instead of $\mathcal{B}_{\mathcal{H}}$. In this case, Proposition 1 is generalized without losing the advantages of the size and time complexities. Detailed complexity analysis is found in Appendix B.

As an application for XML manipulation, we consider the simple example (Fig. 1) that contains the hedge automaton \mathcal{H}_{in} , whose final state is γ . The language accepted by \mathcal{H}_{in} consists of the hedges $f\langle g\langle \mathbf{a}^* \rangle g\langle \mathbf{b}^* \rangle \rangle$. Here \circ is omitted, and \mathbf{a}^* (\mathbf{b}^*) denotes an arbitrary sequence of \mathbf{a} (resp. \mathbf{b}). In the example, we suppose that \mathbf{a} and \mathbf{b} are the elements with different types, and f and g are the tags, such that g contains \mathbf{a} or \mathbf{b} only. According to Lemma 7, we can construct tree automata $\mathcal{A}_{\mathcal{H}_{\text{in}}}$ over $F_{\Sigma, C}$, which corresponds to \mathcal{H}_{in} .

Now we consider the constraint Ψ_{\leq} for hedges, such that $\Psi_{\leq}(t)$ holds if $t = f\langle g\langle a^* \rangle g\langle b^* \rangle \rangle$ and the number of occurrences of a is less than or equals to the number of b . For instance, $\Psi_{\leq}(f\langle g\langle aa \rangle g\langle bbb \rangle \rangle)$ is true.

In tree automata framework, the above predicate Ψ_{\leq} can be translated to the TAN $\mathcal{A}_{\leq} = (\mathcal{R}, Q, Q_{\text{fin}}, \Delta_{\leq})$ with $\mathcal{R}_{\leq} = (F_{\Sigma, C}, \emptyset, R_{\leq})$ as follows.

$$\begin{aligned} \Delta_{\leq}: \quad & f(\eta(x)) \rightarrow \theta(f(x)) \quad \gamma(x) \circ \delta(y) \rightarrow \eta(x \circ y) \\ & g(\alpha(x)) \rightarrow \gamma(g(x)) \quad \epsilon \rightarrow \alpha(\epsilon) \\ & g(\beta(x)) \rightarrow \delta(g(x)) \quad \epsilon \rightarrow \beta(\epsilon) \quad \mathbf{b} \rightarrow \beta(\mathbf{b}) \quad \beta(x) \circ \beta(y) \rightarrow \beta(x \circ y) \\ Q: \quad & \alpha \quad \beta \quad \gamma \quad \delta \quad \eta \quad \theta \\ Q_{\text{fin}}: \quad & \theta \\ R_{\leq}: \quad & f(g(a \circ x) \circ g(b \circ y)) \rightarrow f(g(x) \circ g(y)) \quad f(g(a) \circ g(b)) \rightarrow f(g(\epsilon) \circ g(\epsilon)) \\ & f(g(a \circ x) \circ g(b)) \rightarrow f(g(x) \circ g(\epsilon)) \quad f(g(a) \circ g(b \circ y)) \rightarrow f(g(\epsilon) \circ g(y)) \\ & (x \circ y) \circ z \rightarrow x \circ (y \circ z) \quad \epsilon \circ x \rightarrow x \quad x \circ \epsilon \rightarrow x \end{aligned}$$

Lemma 8. \mathcal{A}_{\leq} with \mathcal{R}_{\leq} accepts t if and only if $\text{flat}(t)$ satisfies Ψ_{\leq} . \square

Using a similar construction, we can define the TAN $\mathcal{A}_{<}$, \mathcal{A}_{\geq} , $\mathcal{A}_{>}$ for the constraints $\Psi_{<}$ (less than), Ψ_{\geq} (greater than or equals to) and $\Psi_{>}$ (greater than).

Remark 3. The languages $\mathcal{L}(\mathcal{A}_{\leq})$, $\mathcal{L}(\mathcal{A}_{<})$, $\mathcal{L}(\mathcal{A}_{\geq})$, $\mathcal{L}(\mathcal{A}_{>})$ are not regular.

Furthermore, $\mathcal{L}(\mathcal{A}_{\leq})$ is closed under Boolean operations and the membership problem is decidable, because the above \mathcal{R}_{\leq} is strongly normalizing and confluent as there are 9 critical pairs and they are all joinable. Besides, the emptiness, universality and inclusion problems are decidable, because R_{\leq} is left-linear and $E = \emptyset$. Similar properties hold for $\mathcal{L}(\mathcal{A}_{<})$, $\mathcal{L}(\mathcal{A}_{\geq})$ and $\mathcal{L}(\mathcal{A}_{>})$.

6 Related Work and Conclusions

Recently, Dal Zilio and Lugiez [7] and thereafter Seidl *et al.* [25] proposed the extension of XML schema language in which a commutative infix operator over hedges is equipped. We denote below this operator by \otimes . The concatenation $t_1 \otimes \cdots \otimes t_n$ of hedges is represented by the multiset $\{t_1, \dots, t_n\}$ in [25]. By this extension, their frameworks require the additional sort of transition rules

$$a\langle \alpha_1(x_1) \otimes \cdots \otimes \alpha_n(x_n) \rangle \rightarrow \beta(a\langle x_1 \otimes \cdots \otimes x_n \rangle) \Leftarrow \alpha_1 \cdots \alpha_n \in \mathcal{L}(\mathcal{G}/C)$$

A *commutative regular grammar* \mathcal{G}/C is a regular grammar modulo the commutativity. The language generated by \mathcal{G}/C is the same as the commutative closure of $\mathcal{L}(\mathcal{G})$, i.e. $u \in \mathcal{L}(\mathcal{G}/C)$ if and only if $u =_C v$ and $v \in \mathcal{L}(\mathcal{G})$. The class of the commutative closure of regular languages coincides with the class of the commutative closure of context-free languages. Moreover, the Parikh image [24] of regular languages (also of the commutative closure of regular languages) is the semi-linear sets, as well as Presburger definable languages. Using this fact,

Dal Zilio and Lugiez define a query language for XML, called *sheaves logic*, in which the formulas may contain Presburger constraints so the number of occurrences of a content with a certain type in a sequence of contents can be inquired.

In tree automata with ranked signature, this extension in hedge automata can be simulated by assuming that \otimes has the associativity $(x \otimes y) \otimes z = x \otimes (y \otimes z)$ and commutativity $x \otimes y = y \otimes x$ axioms. Hedge automata with this \otimes are thus modeled by regular tree automata modulo the AC axioms. However, the example of XML manipulation discussed in Section 5, e.g. Lemma 8, is beyond this class of ETA, and moreover,

Remark 4. the constraints $\Psi_{\leq}, \Psi_{<}, \Psi_{\geq}, \Psi_{>}$ are not expressed by sheaves logic.

Regarding the translation of hedge automata, the special class of tree automata, called *stepwise tree automata*, is studied [3]. The signature of this class consists of constant symbols from Σ and the binary symbol $@$. Their transition rules are in the two forms:

$$a \rightarrow \alpha(a) \qquad \alpha(x) @ \beta(y) \rightarrow \gamma(x @ y)$$

for some $a \in \Sigma$ and α, β, γ are state symbols. Epsilon transition rules $\alpha(x) \rightarrow \beta(x)$ can be permitted though they do not appear formally in the original definition. This class of tree automata has the bidirectional correspondence to hedge automata in terms of Proposition 1. Moreover, we can find in [17] the comprehensive study on the minimization of stepwise automata.

In the research of XML document processing, one of the prominent topics is the *type checking problem*. This question asks us if, for XML document types τ_{in} and τ_{out} and an XML transformer Φ , every document transformed by Φ from of type τ_{in} is of type τ_{out} . In practice, these document types τ_{in} and τ_{out} are described in an XML schema language, and this problem is represented as the inclusion problem $\{t \mid \exists s \in \mathcal{L}(\tau_{\text{in}}): t = \Phi(s)\} \subseteq \mathcal{L}(\tau_{\text{out}})$. For reasoning about this problem, Milo *et al.* [18] studied a special class of transformer on ranked trees, called *k-pebble tree transducer (k-PTT)*, and they showed that the complexity of the type checking problem for *k-PTT* is *k*-tower of exponential. Engelfriet *et al.* showed in [8] that *k-PTT* is simulated by $(k + 1)$ -fold compositions of another kind of transformer, called *stay macro tree transducers (SMTT)*. Recently, Maneth *et al.* showed in [15] that the type checking can be solved in polynomial time if SMTT is *n-bounded copying*, that is, if nodes in an input tree are copied by transformer at most *n*-times in the output tree.

In rewriting, the type checking problem can be reformulated as the problem $\{t \mid \exists s \in \mathcal{L}(\tau_{\text{in}}): s \xrightarrow{!}_{\mathcal{R}_{\Phi}} t\} \subseteq \mathcal{L}(\tau_{\text{out}})$. In [26], Takai *et al.* proposed a class of rewrite systems that effectively preserves the regularity. This result enables us to handle the type checking problem whenever \mathcal{R}_{Φ} is in this class. For automated reasoning, Ohsaki and Takai have developed the verification tool *ACTAS* [21], and Genet *et al.* independently developed another tool *Timbuk* [9], that computes the exact and under-approximated \mathcal{R} -descendant closure of regular tree languages.

In rewriting community, equational rewriting has been studied over the years, and this theory underlies several practical challenges. In fact, there are executable specification languages, verification tools, and program optimization

tools that support rewriting modulo equational axioms. Among elaborated studies on equational rewriting, tree automata with normalization (TAN) may have the roots in *normalized rewriting* [16]. The theory of normalized rewriting brought the idea in dealing with completion modulo axioms, and in practice that developed the software *CiME* [6].

In the paper, we showed the decidability (Theorems 1–4 and Corollaries 1,2) and the closure properties (Theorem 5 and Corollaries 3,4) of our tree automata framework. Corollary 5 can be improved if the intersection-emptiness of AC-regular language and R_1 -normal forms is a decidable question. In Section 5 we discussed advantages of TAN through the observation of transformations from XML schema to tree automata. In the example of XML manipulation, we showed by demonstrating that XML schema and a constraint can be translated to TAN.

For further exploration, we have an interest in the extension of TAN, in the direction of PTA [10] and of monotone ETA [23]. The tool development based on our framework is useful in practice, and it can be easily implemented in the rule-based software, in particular for the tools supporting AC tree automata [21], due to Lemma 1 and Properties (i)–(iii).

Acknowledgments. The authors would like to thank Ralf Treinen and Florent Jacquemard for their comments on the early work. We have done most of the preliminary work while Ohsaki visited in ENS de Cachan in August and September, 2006. We also thank Stéphanie Delaune for her suggestion that improves Lemma 3.

References

1. A. Armando, D. Basin, Y. Boichut, Y. Chevalier, L. Compagna, J. Cuellar, P. Hankes Drielsma, P.-C. Heám, O. Kouchnarenko, J. Mantovani, S. Mödersheim, D. von Oheimb, M. Rusinowitch, J. Santiago, M. Turuani, L. Viganò and L. Vigneron: *The AVISPA Tool for the Automated Validation of Internet Security Protocols and Applications*, Proc. of 17th CAV, Edinburgh (UK), LNCS 3576, pp. 281–285, Springer, 2005.
2. F. Baader and T. Nipkow: *Term Rewriting and All That*, Cambridge University Press, 1998.
3. J. Carme, J. Niehren and M. Tommasi: *Querying Unranked Trees with Stepwise Tree Automata*, Proc. of 15th RTA, Aachen (Germany), LNCS 3091, pp. 105–118, Springer, 2004.
4. H. Comon, M. Dauchet, R. Gilleron, F. Jacquemard, D. Lugiez, S. Tison and M. Tommasi: *Tree Automata Techniques and Applications*, draft, 2005. Available at URL <http://www.grappa.univ-lille3.fr/tata>
5. H. Comon and F. Jacquemard: *Ground Reducibility is EXPTIME-Complete*, Information and Computation 187(1), pp. 123–153, Elsevier, 2003.
6. E. Contejean, C. Marché, B. Monate and X. Urbain: *The CiME System: Version 2.02*, 2004. Software and the document available at URL <http://cime.lri.fr/>
7. S. Dal Zilio and D. Lugiez: *XML Schema, Tree Logic and Sheaves Automata*, Applicable Algebra in Engineering, Communication and Computing 17(5), pp. 337–377, Springer, 2006.
8. J. Engelfriet and S. Maneth: *A Comparison of Pebble Tree Transducers with Macro Tree Transducers*, Acta Informatica 39, pp. 613–698, Springer, 2003.
9. G. Feuillade, T. Genet and V. Viet Triem Tong: *Reachability Analysis over Term Rewriting Systems*, Journal of Automated Reasoning 33, pp. 341–383, Springer, 2004.

10. J. Hendrix, H. Ohsaki and M. Viswanathan: *Propositional Tree Automata*, Proc. of 17th RTA, Seattle (Washington), LNCS 4098, pp. 50–65, Springer, 2006.
11. H. Hosoya, J. Vouillon and B.C. Pierce: *Regular Expression Types for XML*, Proc. of 5th ICFP, Montreal (Canada), SIGPLAN Notices 35(9), pp. 11–22, ACM Press, 2000.
12. N. Klarlund and A. Møller: *MONA Version 1.4 User Manual*, BRICS Notes Series NS-01-1, Department of Computer Science, University of Aarhus, 2001.
13. T. Kutsia: *Solving Equations Involving Sequence Variables and Sequence Functions*, Proc. of 7th AISC, Linz (Austria), LNCS 3249, pp. 157–170, Springer, 2004.
14. D. Lugiez: *Multitree Automata That Count*, Theoretical Computer Science 333(1–2), pp. 225–263, Elsevier, 2005.
15. S. Maneth, T. Perst and H. Seidl: *XML Type Checking in Polynomial Time*, Proc. of 11th ICDT, Barcelona (Spain), LNCS 4353, pp. 254–268, Springer, 2007.
16. C. Marché: *Normalized Rewriting: An Alternative to Rewriting Modulo a Set of Equations*, Journal of Symbolic Computation 21(3), pp. 253–288, Elsevier, 1996.
17. W. Martens and J. Niehren: *Minimizing Tree Automata for Unranked Trees*, Proc. of 10th DBPL, Trondheim (Norway), LNCS 3774, pp. 232–246, Springer, 2005.
18. T. Milo, D. Suciú and V. Vianu: *Typechecking for XML Transformers*, Proc. of 19th PODS, Dallas (Texas), pp. 11–22, ACM Press, 2000.
19. M. Murata: *Extended Path Expression for XML*, Proc. of 20th PODS, Santa Barbara (California), pp. 61–72, ACM Press, 2001.
20. M. Murata and P. Prescod: *SGML/XML and Forest/Hedge Automata Theory*, Online Resource for Markup Language Technologies, Cover Pages, OASIS, May 2001. Document available at URL <http://xml.coverpages.org/hedgeAutomata.html>
21. H. Ohsaki and T. Takai: *ACTAS : A System Design for Associative and Commutative Tree Automata Theory*, Proc. of 5th RULE, Aachen (Germany), ENTCS 124, pp. 97–111, Elsevier, 2005.
22. H. Ohsaki and T. Takai: *Decidability and Closure Properties of Equational Tree Languages*, Proc. of 13th RTA, Copenhagen (Denmark), LNCS 2378, pp. 114–128, Springer, 2002.
23. H. Ohsaki: *Beyond Regularity: Equational Tree Automata for Associative and Commutative Theories*, Proc. of 15th CSL, Paris (France), LNCS 2142, pp. 539–553, Springer, 2001.
24. R. Parikh: *On Context-Free Languages*, JACM 13(4), pp. 570–581, ACM Press, 1966.
25. H. Seidl, T. Schwentick and A. Muscholl: *Numerical Document Queries*, Proc. of 22nd PODS, San Diego (California), pp. 155–166, ACM Press, 2003.
26. T. Takai, Y. Kaji and H. Seki: *Right-Linear Finite Path Overlapping Term Rewriting Systems Effectively Preserve Recognizability*, Proc. of 11th RTA, Norwich (UK), LNCS 1833, pp. 246–260, Springer, 2000.
27. J.W. Thatcher: *Characterizing Derivation Trees of Context-Free Grammars Through a Generalization of Automata Theory*, Journal of Computer and System Sciences 1(4), pp. 317–322, Elsevier, 1967.
28. K.N. Verma: *Two-Way Equational Tree Automata for AC-Like Theories: Decidability and Closure Properties*, Proc. of 14th RTA, Valencia (Spain), LNCS 2706, pp. 180–197, Springer, 2003.
29. K.N. Verma: *On Closure under Complementation of Equational Tree Automata for Theories Extending AC*, Proc. of 10th LPAR, Almaty (Kazakhstan), LNCS 2850, pp. 183–197, Springer, 2003.
30. *Extensible Markup Language (XML) Schema 1.0*, 3rd edition, W3C, 2004. Document available at URL <http://www.w3.org/TR/2004/REC-xml-20040204/>

A Proofs

Proof of Remark 1. For the class of equational tree automata, Lemma 9 hold.

Lemma 9. *If $R = \emptyset$, $\mathcal{L}(\mathcal{A}) = \mathcal{E}(\mathcal{L}(\mathcal{A}'))$ where \mathcal{A}' is obtained by replacing E in \mathcal{A} with \emptyset .*

One should note that this does not hold if we take the following syntax

$$f(\alpha_1, \dots, \alpha_n) \rightarrow \beta \quad \alpha \rightarrow \beta$$

as the definition of transitions rules.

We recall the proof by Verma [29]. Suppose for leading to the contradiction that there exists an automaton \mathcal{B} with the above R and E in Remark 1, such that \mathcal{B} accepts $\mathcal{L}(\mathcal{A}_2)$. Let \mathcal{B}' be the tree automaton obtained by replacing E in \mathcal{B} with \emptyset , then by the above lemma, $\mathcal{L}(\mathcal{B}) = \mathcal{E}(\mathcal{L}(\mathcal{B}'))$. This implies that for every i, j with $i \neq j$, there exists a term t which is accepted by \mathcal{B}' and t is equivalent to

$$f[\mathbf{h}(\mathbf{g}^i(\mathbf{a})), \dots, \mathbf{h}(\mathbf{g}^i(\mathbf{a})), \mathbf{h}(\mathbf{g}^j(\mathbf{a})), \dots, \mathbf{h}(\mathbf{g}^j(\mathbf{a}))]$$

modulo $\text{AC}(f)$. Here $f[\]$ represents a context consisting only of f . From PIGEON-HOLE PRINCIPLE, there exist i and j such that $i \neq j$ and $\{\alpha \in Q \mid \mathbf{g}^i(\mathbf{a}) \xrightarrow{*}_{\mathcal{B}'} \alpha(\mathbf{g}^i(\mathbf{a}))\} = \{\alpha \in Q \mid \mathbf{g}^j(\mathbf{a}) \xrightarrow{*}_{\mathcal{B}'} \alpha(\mathbf{g}^j(\mathbf{a}))\}$. Hence, the above term is accepted by \mathcal{B}' if and only if $f[\mathbf{h}(\mathbf{g}^i(\mathbf{a})), \dots, \mathbf{h}(\mathbf{g}^i(\mathbf{a})), \mathbf{h}(\mathbf{g}^i(\mathbf{a})), \dots, \mathbf{h}(\mathbf{g}^i(\mathbf{a}))]$ is accepted. \square

Proof of Lemma 1. First we define the mapping rem for terms in $T(F \cup Q)$ such that

$$\text{rem}(t) = \begin{cases} f(\text{rem}(t_1), \dots, \text{rem}(t_n)) & \text{if } t = f(t_1, \dots, t_n) \text{ for some } f \in F, \\ \text{rem}(t') & \text{if } t = \alpha(t') \text{ for some } \alpha \in Q. \end{cases}$$

This mapping strips off all state symbols from the input. By structural induction, we can prove that $\text{rem}(t) \in \text{NF}_F(\mathcal{R})$ implies $t \in \text{NF}_F(\mathcal{R})$. Next we prove this lemma and additionally the property that $\text{rem}(t') =_{\mathcal{E}} (t) \downarrow_{\mathcal{R}}^!$. We use induction on the length of $\rightarrow_{\mathcal{A}}^+$. The base case is easy, because $\rightarrow_{\mathcal{A}}$ and $\rightarrow_{\mathcal{R}}^! \rightarrow_{\mathcal{A}'}$ coincide by definition. Moreover, in this case, $(t) \downarrow_{\mathcal{R}}^!$ and t' can be represented as $(t) \downarrow_{\mathcal{R}}^! =_{\mathcal{E}} C[c]$ and $t' =_{\mathcal{E}} C[\alpha(c)]$ for some context C over F , constant c from F and state α in Q . If $t' \notin \text{NF}_F(\mathcal{R})$, then $C[\alpha(c)] \notin \text{NF}_F(\mathcal{R})$, that implies $C[c] \notin \text{NF}_F(\mathcal{R})$, but it contradicts to the assumption. For the induction step, we assume $t \rightarrow_{\mathcal{A}}^+ t''$. By induction hypothesis, $(t) \downarrow_{\mathcal{R}}^! \xrightarrow{*}_{\mathcal{A}'} t''$ and $\text{rem}(t'') =_{\mathcal{E}} (t) \downarrow_{\mathcal{R}}^!$, and then $t'' \in \text{NF}_F(\mathcal{R})$. Thus, $(t) \downarrow_{\mathcal{R}}^! \xrightarrow{*}_{\mathcal{A}'} t'' \rightarrow_{\mathcal{A}'} t'$. Regarding the transition $t'' \rightarrow_{\mathcal{A}'} t'$, we have the two cases:

- (1) $t'' =_{\mathcal{E}} C'[f(\alpha_1(u_1), \dots, \alpha_n(u_n))]$ and $t' =_{\mathcal{E}} C'[\beta(f(u_1, \dots, u_n))]$ for some $f(\alpha_1(x_1), \dots, \alpha_n(x_n)) \rightarrow \beta(f(x_1, \dots, x_n)) \in \Delta$,
- (2) $t'' =_{\mathcal{E}} C'[\alpha(u)]$ and $t' =_{\mathcal{E}} C'[\beta(u)]$ for some $\alpha(x) \rightarrow \beta(x) \in \Delta$.

Here C' is a context over $F \cup Q$, u_1, \dots, u_n, u are terms over $F \cup Q$, and $\alpha_1, \dots, \alpha_n, \alpha, \beta$ are states in Q . In any case, by induction hypothesis, since $\text{rem}(t'') =_{\mathcal{E}} (t) \downarrow_{\mathcal{R}}^!$, we have $\text{rem}(t') =_{\mathcal{E}} (t) \downarrow_{\mathcal{R}}^!$, because $\text{rem}(t'') =_{\mathcal{E}} \text{rem}(t')$. \square

Proof of Theorem 3. Given tree automata \mathcal{A}_1 and \mathcal{A}_2 with \mathcal{R} -normalization. From assumption, we have an equational tree automaton \mathcal{B} with the same equational theory which accepts $\text{NF}_F(\mathcal{R})$. From monotonicity of closure operators, we obtain

$$\mathcal{L}(\mathcal{A}_1) \subseteq \mathcal{L}(\mathcal{A}_2) \Leftrightarrow \mathcal{L}(\mathcal{A}'_1) \cap \mathcal{L}(\mathcal{B}) \subseteq \mathcal{L}(\mathcal{A}'_2) \cap \mathcal{L}(\mathcal{B})$$

where \mathcal{A}'_i for $i \in \{1, 2\}$ is an equational tree automaton define by replacing R in \mathcal{A}_i with \emptyset \square

Proof of Lemma 2. Let $\mathcal{A} = (\mathcal{R}, Q, Q_{\text{fin}}, \Delta)$. From assumption, $(t) \downarrow_{\mathcal{R}}^!$ is unique up to modulo \mathcal{E} . We take $\mathcal{A}' = ((\mathcal{E}, \emptyset), Q, Q_{\text{fin}}, \Delta)$. From Lemma 1, we have

$$t \in \mathcal{L}(\mathcal{A}) \Leftrightarrow (t) \downarrow_{\mathcal{R}}^! \in \mathcal{L}(\mathcal{A}') \Leftrightarrow (t) \downarrow_{\mathcal{R}}^! \notin (\mathcal{L}(\mathcal{A}'))^c.$$

The second equivalence is obtained from the confluence of \mathcal{R} . From assumption, there effectively exists an equational tree automaton $\mathcal{B} = ((\mathcal{E}, \emptyset), Q', Q'_{\text{fin}}, \Delta')$ such that $\mathcal{L}(\mathcal{B}) = (\mathcal{L}(\mathcal{A}'))^c$. Therefore, $((\mathcal{E}, R), Q', Q'_{\text{fin}}, \Delta')$ accepts t if and only if \mathcal{A} does not accept t . \square

Proof: Why strong normalization for Theorem 1? We use the typical reduction from PCP (Post's correspondence problem). Given an instance of the PCP $P = \{(u_i, v_i) \mid 1 \leq i \leq n: u_i, v_i \in \{a, b\}^+\}$, the set R_P consists of the rewrite rules

$$f(x, y) \rightarrow f(\overline{u_i}(x), \overline{v_i}(y)) \quad \text{for each } (u_i, v_i) \in P$$

and

$$f(a(x), a(y)) \rightarrow g(a(x), a(y)) \quad f(b(x), b(y)) \rightarrow g(b(x), b(y))$$

$$g(a(x), a(y)) \rightarrow g(x, y) \quad g(b(x), b(y)) \rightarrow g(x, y)$$

$$g(c, c) \rightarrow d \quad f(x, y) \rightarrow e$$

where $\overline{u}(t) = \overline{u'}(a(t))$ if $u = a u'$ and $\overline{u}(t) = \overline{u'}(b(t))$ if $u = b u'$; otherwise, $\overline{u}(t) = t$. Let $\mathcal{R} = ((\{a, b, c, d, e, f, g\}, \emptyset), R_P)$ and $\mathcal{A}_P = (\mathcal{R}, \{\alpha\}, \{\alpha\}, \{d \rightarrow \alpha(d)\})$. Observe that \mathcal{R} is weakly normalizing, because every term of the form $f(t_1, t_2)$ can be rewritten to e , and \mathcal{R} is strongly normalizing over $\mathcal{T}(F - \{f\})$. From Lemma 1, we can see that $f(c, c) \in \mathcal{L}(\mathcal{A}_P)$ if and only if $(f(c, c)) \downarrow_{\mathcal{R}}^! = d$. Trivially, the latter statement holds if and only if P has a solution. \square

Proof of Theorem 5. It is obvious that this class of tree languages is closed under union. For the intersection, we take an equational rewriting system $\mathcal{R} = (\mathcal{E}, R)$ with $\mathcal{E} = (F, E)$, and tree automata \mathcal{A} and \mathcal{B} modulo \mathcal{R} -normalization. Define \mathcal{A}' and \mathcal{B}' to be equational tree automata, each of which is obtained by replacing R with \emptyset , respectively for \mathcal{A} and \mathcal{B} . Then, from assumption, there exists $\mathcal{C} = ((\mathcal{E}, \emptyset), Q, Q_{\text{fin}}, \Delta)$ such that $\mathcal{L}(\mathcal{C}) = \mathcal{L}(\mathcal{A}') \cap \mathcal{L}(\mathcal{B}')$. For every $t \in \mathcal{T}(F)$, $(t) \downarrow_{\mathcal{R}}^!$ is unique, so $(t) \downarrow_{\mathcal{R}}^! \in \mathcal{L}(\mathcal{A}')$ and $(t) \downarrow_{\mathcal{R}}^! \in \mathcal{L}(\mathcal{B}')$ if and only if $t \in \mathcal{L}(\mathcal{A}) \cap \mathcal{L}(\mathcal{B})$. Thus, $((\mathcal{E}, R), Q, Q_{\text{fin}}, \Delta)$ accepts $\mathcal{L}(\mathcal{A}) \cap \mathcal{L}(\mathcal{B})$. \square

Proof of Lemma 6. Without loss of generality we assume that every production rule in \mathcal{A} is in the shape of $\alpha \rightarrow a\beta$ or $\alpha \rightarrow a$. Hereafter in the proof, we show

the tree automata construction over the signature $F = \Sigma \cup \{f\}$ with $\text{arity}(f) = 2$. We define the tree automaton $\mathcal{A} = (F, Q, Q_{\text{fin}}, \Delta)$ associated to \mathcal{G} as follows:

$$\begin{aligned} F &: \Sigma \cup \{f\} \text{ such that } \text{arity}(a) = 0 \text{ for all } a \in \Sigma \\ Q &: \{q_{\alpha \rightarrow \beta} \mid \alpha, \beta \in \mathcal{S}\} \cup \{q_\alpha \mid \alpha \in \mathcal{S}\} \\ Q_{\text{fin}} &: \{q_{s_0}\} \\ \Delta &: \{f(q_{\alpha \rightarrow \beta}(x), q_\beta(y)) \rightarrow q_\alpha(f(x, y)) \mid \alpha, \beta \in \mathcal{S}\} \cup \\ &\quad \{f(q_{\alpha \rightarrow \beta}(x), q_{\beta \rightarrow \gamma}(y)) \rightarrow q_{\alpha \rightarrow \gamma}(f(x, y)) \mid \alpha, \beta, \gamma \in \mathcal{S}\} \cup \\ &\quad \{a \rightarrow q_{\alpha \rightarrow \beta}(a) \mid \alpha \rightarrow a\beta \in \Lambda \text{ for some } \alpha, \beta \in \mathcal{S} \text{ and } a \in \Sigma\} \cup \\ &\quad \{a \rightarrow q_\alpha(a) \mid \alpha \rightarrow a \in \Lambda \text{ for some } \alpha \in \mathcal{S} \text{ and } a \in \Sigma\} \end{aligned}$$

Then we can show by structural induction that t is accepted by \mathcal{A} if and only if $\text{leaf}(t)$ is generated by \mathcal{G} . By construction of Δ , we observe that for all $w \in \Sigma^*$, $a \in \Sigma$ and $\alpha, \beta \in \mathcal{S}$, (1) $\alpha \rightarrow_{\mathcal{G}} a$ if and only if $a \rightarrow_{\mathcal{A}} q_\alpha(a)$, and (2) $\alpha \rightarrow_{\mathcal{G}}^* w\beta$ if and only if $u \rightarrow_{\mathcal{A}}^* q_{\alpha \rightarrow \beta}(u)$ and $\text{leaf}(u) = w$. \square

Proof: Why Lemma 6 does not hold if \mathcal{G} is context-free grammar? We consider the language $L = \{a^n b^n \mid n \geq 1\}$ which is generated by a context-free grammar. We suppose for leading the contradiction that there is a tree automaton $\mathcal{A} = (F, Q, Q_{\text{fin}}, \Delta)$ whose leaf language is L . The signature F must contain a function symbol whose arity is greater than 1. But since F is finite, trees accepted by \mathcal{A} do not have an upper-bound in the depth. In fact, when we take a natural number k which is *large* enough to apply PUMPING LEMMA, there exists a tree $C[t]$ whose subtree t satisfies that $\text{depth}(t) > |Q|$, $\text{leaf}(t) = b^+$ and $\text{length}(\text{leaf}(t)) > |Q|$. Then \mathcal{A} accepts $C[t']$ whose subtree t' has less leaves than that of t . Because $\text{leaf}(C[t'])$ is no longer a word in L , it leads to the contradiction. \square

Proof of Lemma 7. We define the grammar $\mathcal{G} = (\Sigma, s_0, \mathcal{S}, \Lambda)$ associated to \mathcal{A} as follows:

$$\begin{aligned} \Sigma &: \{f \in F \mid \text{arity}(f) = 0\} \\ \mathcal{S} &: Q \cup \{s_0\} \\ \Lambda &: \left\{ \beta \rightarrow c_1 \dots c_n \alpha \mid \begin{array}{l} f(c_1, \dots, c_n, \alpha(x)) \rightarrow_{\mathcal{A}}^* \beta(f(c_1, \dots, c_n, x)) \\ \text{for some } f, c_1, \dots, c_n \in F \text{ and } \alpha, \beta \in Q \end{array} \right\} \cup \\ &\quad \{ \beta \rightarrow c \mid c \rightarrow_{\mathcal{A}}^* \beta(c) \text{ for some } c \in F \text{ and } \beta \in Q \} \cup \\ &\quad \{ s_0 \rightarrow \alpha \mid \alpha \in Q_{\text{fin}} \} \end{aligned}$$

Using the induction on t in $\text{NF}_F(\mathcal{R}_{\text{reg}})$, we can prove that t is accepted by \mathcal{A} if and only if $\text{leaf}(t)$ is generated by \mathcal{G} . Since \mathcal{G} generates a regular language, due to Corollary 6, we can construct a tree automaton \mathcal{B} that accept trees whose leaves are generated by \mathcal{G} . Hence \mathcal{B} accepts $\{t \mid (t) \downarrow_{\mathcal{R}_{\text{reg}}}^! \text{ accepted by } \mathcal{A}\}$. \square

B Complexity

For an equational rewrite system $\mathcal{R} = (F, E, R)$, if $E = \emptyset$ and $R = \emptyset$, we write $\mathcal{A} = (F, Q, Q_{\text{fin}}, \Delta)$ for $\mathcal{A} = ((F, \emptyset, \emptyset), Q, Q_{\text{fin}}, \Delta)$, and we say \mathcal{A} is a tree

automaton. Let us estimate the size of the tree automata and word grammar constructed in the proofs, and evaluate the computational complexity needed for the constructions.

In the proof of Lemma 7, an input tree automaton \mathcal{A} is translated into a regular word grammar \mathcal{G} , and then a tree automaton \mathcal{B} is obtained due to Corollary 6. In Proposition 1, we explain the property that a hedge automaton \mathcal{H} can be translated into a tree automaton \mathcal{A} , and this statement is generalized in Corollary 7, as a hedge automaton \mathcal{H} is translated into a tree automaton \mathcal{B} via \mathcal{A} . Hereafter in the following, we let $\mathcal{H} = (\Sigma, C, N, N_{\text{fin}}, \Delta_{\mathcal{H}})$, $\mathcal{A} = (F, Q, Q_{\text{fin}}, \Delta)$, $\mathcal{B} = (F', Q', Q'_{\text{fin}}, \Delta')$, and $\mathcal{G} = (\Sigma, \mathcal{S}, s_0, A)$.

Complexity in Corollary 6. In the proof of Lemma 6, the number of constructed rules is $O(|A| + |\mathcal{S}|^3)$ where the factor $|\mathcal{S}|^3$ is due to the second line of the constructed rules Δ in the proof. When Lemma 6 is extended to Corollary 6, the number of the rules grows up to $O(|A| + |\mathcal{S}|^{N+1})$ where N is the maximum arity of function symbols in F . Precisely, $|F'| \in O(|\Sigma|)$, $|Q'| \in O(|\mathcal{S}|^2)$, $|Q'_{\text{fin}}| \in O(1)$ and $|\Delta'| \in O(|A| + |\mathcal{S}|^{N+1})$.

Complexity in Lemma 7. Let us look at the constructed rules A in the proof of Lemma 7. The number of rules in the first line of A is $O(|F|^{N-1} \times |Q|^2)$. This exponential blowing-up can be avoided by constructing a rule of the shape $\beta \rightarrow c\gamma$ ($\beta, \gamma \in Q$, $c \in F$) one by one instead of constructing a rule of the shape $\beta \rightarrow c_1 \dots c_n \alpha$. Thus, $|A| \in O(N|F||\Delta| + |Q_{\text{fin}}|)$. Combining the size estimation in Corollary 6 and this observation, the size of \mathcal{B} can be represented by $|F'| \in O(|\Sigma|) = O(|F|)$, $|Q'| \in O(|\mathcal{S}|^2) = O(|Q|^2)$, $|Q'_{\text{fin}}| \in O(1)$ and $|\Delta'| \in O(|A| + |\mathcal{S}|^{N+1}) = O(N \times |F| \times |\Delta| + |Q_{\text{fin}}| + |Q|^{N+1})$. The time complexity for the construction of \mathcal{B} from \mathcal{A} is also $O(N \times |F| \times |\Delta| + |Q_{\text{fin}}| + |Q|^{N+1})$.

Complexity in Proposition 1. The construction of $|\mathcal{A}|$ in Proposition 1 is rather straightforward and $|F| \in O(|\Sigma| + |C|)$, $|Q| \in O(|N| + \sum_{\mathcal{G} \in \Delta_{\mathcal{H}}} |N_{\mathcal{G}}|)$, $|Q_{\text{fin}}| \in O(|N_{\text{fin}}|)$ and $|\Delta| \in O(|\Delta_{\mathcal{H}}| + \sum_{\mathcal{G} \in \Delta_{\mathcal{H}}} |A_{\mathcal{G}}|)$ where $N_{\mathcal{G}}$ (resp. $A_{\mathcal{G}}$) denotes the set of non-terminals (resp. rules) of a regular word grammar \mathcal{G} appearing in a rule in $\Delta_{\mathcal{H}}$. Thus, $|\mathcal{A}| \in O(\mathcal{H})$ and the time needed for the construction of \mathcal{A} is also linear in the size of \mathcal{H} .

Complexity in Corollary 7. The size of each component \mathcal{B} can be estimated by those for Proposition 1 and Lemma 7 above. Note that since the signature $F_{\Sigma, C}$ of \mathcal{A} constructed in the proof of Proposition 1 contains only a binary function symbol and constants, we can let $N = 2$ when we apply Lemma 7 to \mathcal{A} to obtain a desirable tree automaton \mathcal{B} . Therefore, $|F'| \in O(|\Sigma| + |C|)$, $|Q'| \in O((|N| + \sum_{\mathcal{G} \in \Delta_{\mathcal{H}}} |N_{\mathcal{G}}|)^2)$, $|Q'_{\text{fin}}| \in O(1)$, and $|\Delta'| \in O((|\Sigma| + |C|) \times (|\Delta_{\mathcal{H}}| + \sum_{\mathcal{G} \in \Delta_{\mathcal{H}}} |A_{\mathcal{G}}|) + |N_{\text{fin}}| + (|N| + \sum_{\mathcal{G} \in \Delta_{\mathcal{H}}} |N_{\mathcal{G}}|)^3)$. Summarizing, $|\mathcal{B}| \in O(|\mathcal{H}|^3)$ and the time complexity needed for the construction of \mathcal{B} is also $O(|\mathcal{H}|^3)$.