

Normal Form on Linear Tree-to-word Transducers

Adrien Boiret^{1,2}

¹ University Lille 1, France

² Links (Inria Lille & CRISAL, UMR CNRS 9189), France

Abstract. We study a subclass of tree-to-word transducers: linear tree-to-word transducers, that cannot use several copies of the input. We aim to study the equivalence problem on this class, by using minimization and normalization techniques. We identify a Myhill-Nerode characterization. It provides a minimal normal form on our class, computable in EXPTIME. This paper extends an already existing result on tree-to-word transducers without copy or reordering (sequential tree-to-word transducers), by accounting for all the possible reorderings in the output.

1 Introduction

Transducers and their properties have long been studied in various domains of computer sciences. The views on transducers that motivate this paper's field of research are mostly the result of the intersection of two approaches.

Language theory sees transducers as the natural extension of automata, with an output. This view extends almost as far back as the study of regular languages, and developed techniques to solve classical problems such as equivalence, type-checking, or even learning problems (e.g. [11,10,4]) on increasingly wide classes of transducers.

Functional programming sees transducers as a formal representation of some programs. In order to study languages such as XSLT, XQuery, or XProc, used to transform XML trees, classes of transducers that acted more and more like functional programs were designed and studied. For example, deterministic top-down tree transducers can be seen as a functional program that transform trees from the root to the leaves, with finite memory. Different classes extend the reach of transducers to encompass more of the functionalities of programming languages.

Concatenation in the output, notably, plays an important role in the way XSLT produces its outputs. Classes like macro-tree transducers [5], tree-to-word transducers, or even word-to-word transducers with copies in the output [1] allow such concatenation, but as this functionality appears to be difficult to combine with the classical techniques of language theory, this is to the cost of very few results carrying to these classes.

Tree-to-word transducers and Macro-tree transducers are of particular relevance, as they allow concatenation in their output, and are at the current

frontier between the language theory approach of transducers and the approach of transducers seen as functional programs.

Many problems are left open in these classes. Notably, in the general case for Macro-tree transducers, the decidability equivalence is a famous long-standing question, that has yet to be resolved. However, some pre-existing results exist for fragments of these classes.

Equivalence for the subclass of linear size increase macro-tree transducers [3] is proven to be decidable. It comes from a *logic* characterization, as if we bound the number of times a transducer can copy the same subtree in the output, then we limit the expressivity of macro-tree transducers into MSO-definable translations, where equivalence is decidable in non-elementary complexity [4].

Equivalence for all tree-to-word transducers has recently been proven to be decidable in randomized polynomial time [13]. Note that this result uses neither classic logic methods nor the classic transducer methods, and does not provide a characterization or Myhill-Nerode theorem.

Equivalence is PTIME for sequential tree-to-word transducers [6], that prevents copying in the output and forces subtrees to produce following the order of the input. Furthermore, using a *Myhill-Nerode* characterization, a normal form computable in EXPTIME is shown to exist. This normal form was later proven to be learnable in PTIME [7].

In this paper, we aim to study the linear tree-to-word transducers (or LTWs), a restriction of deterministic tree-to-word transducers that forbids copying in the output, but allows the image of subtrees to be flipped in any order. This is a more general class than sequential tree-to-word transducers, but still less descriptive than general tree-to-word transductions. In this class, we show the existence of a normal form, computable in EXPTIME.

Note that even if equivalence is already known to be decidable in a reasonable complexity, finding a normal form is of general interest in and of itself. For example, in [11,8,7], normal forms on transducers defined using a Myhill-Nerode theorem are used to obtain a learning algorithm.

To define a normal form on LTWs, we start by the methods used for sequential tree-to-words transducers (STWs) in [6]. We consider the notion of *earliest* STWs, which normalizes the output production. We can extend this notion to LTWs and study only earliest LTWs without any loss of expressivity.

In [6], this is enough to obtain a Myhill-Nerode characterization. However, by adding the possibility to flip subtree images to LTWs, we created another way for equivalent transducers to differ. The challenge presented by the extension of the methods of [6] becomes to resolve this new degree of freedom, in order to obtain a good normal form with a Myhill-Nerode characterization.

Outline. After introducing basic notions on words and trees, we will present our class of linear tree-to-word transducers in Section 2. Then in Section 3 we will extend the notion of *earliest* production in [6] to the linear case, and find out that we can also extend the algorithm that takes a transducer and compute and equivalent earliest one. However, this is no longer sufficient, as transducers can now also differ in the order they produce their subtrees' output in. Section 4

will detail exactly how two earliest transducers can still differ, by categorizing all possible flips. Finally, Section 5 will compile these results into a Myhill-Nerode theorem. This will allow us to establish a normal form, computable in EXPTIME. We will conclude by a brief recap of the result, and propose several possible next steps for this line of research.

2 Preliminaries

Words and Trees

We begin by fixing notations on standard notions over words and ranked trees.

Words. For a finite set of symbols Δ , we denote by Δ^* the set of finite words over Δ with the concatenation operator \cdot and the empty word ε . For a word u , $|u|$ is its length. For a set of words L , we denote $\text{lcp}(L)$ the longest word u that is a prefix of every word in L , or *largest common prefix*. Also, $\text{lcs}(L)$ is the largest common suffix of L . For $w = u \cdot v$, the left quotient of w by u is $u^{-1} \cdot w = v$, and the right quotient of w by v is $w \cdot v^{-1} = u$.

Context-Free Grammars. For a word alphabet Δ , a context-free grammar (or CFG) is a tuple $G = \{N, S, R\}$ where

- N is a finite set of non-terminals.
- $S \in N$ is the initial non-terminal.
- R is a finite set of rules of form $A \rightarrow u_0 A_1 \dots A_n u_n$, where u_0, \dots, u_n are words of Δ^* , and A_1, \dots, A_n are non-terminals of N .

The languages L_A described by the non-terminal A is recursively defined as the smallest language such that if $A \rightarrow u_0 A_1 \dots A_n u_n \in R$, and $w_1 \in L_{A_1}, \dots, w_n \in L_{A_n}$, then $u_0 w_1 \dots w_n u_n \in L_A$.

We say that G describes $L_G = L_S$.

Ranked Trees. A *ranked alphabet* is a finite set of ranked symbols $\Sigma = \bigcup_{k \geq 0} \Sigma^{(k)}$, where $\Sigma^{(k)}$ is the set of k -ary symbols. Every symbol has a unique arity. A *tree* is a ranked ordered term over Σ . For example, $t = f(a, g(b))$ is a tree over Σ if $f \in \Sigma^{(2)}$, $g \in \Sigma^{(1)}$, $a, b \in \Sigma^{(0)}$. The set of all trees on Σ is T_Σ .

Linear Tree-to-Word Transducers

We define linear tree-to-word transducer, that define a function from T_Σ to Δ^* .

Definition 1. A *linear tree-to-word transducer* (LTW) is a tuple $M = \{\Sigma, \Delta, Q, ax, \delta\}$ where

- Σ is a tree alphabet,
- Δ is a finite word alphabet of output symbols,
- Q is a finite set of states,
- ax is a axiom of form $u_0 q u_1$, where $u_0, u_1 \in \Delta^*$ and $q \in Q$,

– δ is a set of rules of the form

$$q, f \rightarrow u_0 q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)}) u_n$$

where $q, q_1, \dots, q_n \in Q$, $f \in \Sigma$ of rank n and $u_0 \dots u_n \in \Delta^*$; σ is a permutation on $\{1, \dots, n\}$. There is at most one rule per pair q, f .

We define recursively the function $\llbracket M \rrbracket_q$ of a state q . $\llbracket M \rrbracket_q(f(t_1 \dots t_n))$ is

- $u_0 \llbracket M \rrbracket_{q_1}(t_{\sigma(1)}) u_1 \dots \llbracket M \rrbracket_{q_n}(t_{\sigma(n)}) u_n$,
- if $q, f \rightarrow u_0 q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)}) u_n \in \delta$
- undefined, if there is no rule for q, f in δ .

The function $\llbracket M \rrbracket$ of a transducer M with axiom $u_0 q u_1$ is defined as

$$\llbracket M \rrbracket(s) = u_0 \llbracket M \rrbracket_q(s) u_1.$$

Note that to get the definition of STWs as made in [6], we just have to impose that in every rule, σ is the identity.

Example 2. Consider the function $\llbracket M \rrbracket : t \mapsto 0^{|t|}$, that counts the number of nodes in t and writes a 0 in the output for each of them. Our LTW has only one state q , and its axiom is $ax = q$

$$\begin{aligned} q(f(x_1, x_2)) &\rightarrow 0 \cdot q(x_1) \cdot q(x_2) \\ q(a) &\rightarrow 0, \quad q(b) \rightarrow 0 \end{aligned}$$

The image of $f(a, b)$ is $\llbracket M \rrbracket(f(a, b)) = \llbracket M \rrbracket_q(f(a, b))$, using the axiom. Then we use the first rule to get $0 \cdot \llbracket M \rrbracket_q(a) \cdot \llbracket M \rrbracket_q(b)$, and finally, $0 \cdot 0 \cdot 0$

We denote with $\text{dom}(\llbracket M \rrbracket)$ the domain of a transducer M , i.e. all trees such that $\llbracket M \rrbracket(t)$ is defined. Similarly, $\text{dom}(\llbracket M \rrbracket_q)$ is the domain of state q . We note L_q the production of q , i.e. the set $\{\llbracket M \rrbracket_q(t) \mid t \in \text{dom}(\llbracket M \rrbracket_q)\}$.

We define accessibility between states as the transitive closure of appearance in a rule. This means q is accessible from itself, and if there is a rule $q, f \rightarrow u_0 q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)}) u_n$, and q accessible from q' , then all states q_i , $1 \leq i \leq n$, are accessible from q' .

A first remark is that state productions are languages of some CFG:

Lemma 3. *One can compute in polynomial time a CFG that describes L_q .*

Proof. For a LTW $M = \{\Sigma, \Delta, Q, \text{ax}, \delta\}$, we create a grammar $G = \{N, S_q, R\}$ where

- For each $q' \in Q$ there is a matching $S_{q'} \in N$.
- For each rule $q', f \rightarrow u'_0 q'_1(x_{\sigma'(1)}) \dots q'_n(x_{\sigma'(n)}) u'_n$, there is a matching rule $S_{q'} \rightarrow u'_0 S_{q'_1} \dots S_{q'_n} u'_n$

We show that for all $S_{q'}$, $L_{S_{q'}} = L_{q'}$ by induction on the derivation in the CFG. $w \in S_{q'}$ if and only if there exists $S_{q'} \rightarrow u'_0 S_{q'_1} \dots S_{q'_n} u'_n \in R$ and $w_1 \in L_{S_{q'_1}}, \dots, w_n \in L_{S_{q'_n}}$ such that $w = u'_0 w_1 \dots w_n u'_n$. By induction, $w_1 \in L_{q'_1}, \dots, w_n \in L_{q'_n}$. This means that $w \in S_{q'}$ if and only if there exists a rule $q', f \rightarrow u'_0 q'_1(x_{\sigma'(1)}) \dots q'_n(x_{\sigma'(n)}) u'_n \in \delta$ and trees $t_{\sigma'(1)}, \dots, t_{\sigma'(n)}$ such that $w = u'_0 \llbracket M \rrbracket_{q'_1}(t_{\sigma'(1)}) \dots \llbracket M \rrbracket_{q'_n}(t_{\sigma'(n)}) u'_n$. Hence $w \in S_{q'}$ if and only if there exists a tree $t = f(t_1, \dots, t_n)$ such that $\llbracket M \rrbracket_{q'}(t) = w$.

We start the normalization process with a natural notion of trimmed LTWs.

Definition 4. A LTW is trimmed if its axiom is $u_0q_0v_0$, and every state q is accessible from q_0 and of non-empty domain.

Note that all LTWs can be made trimmed by deleting all their useless states.

Lemma 5. For M a LTW, one can compute an equivalent trimmed LTW in linear time.

3 Earliest Linear Transducers

It is possible for different LTWs to encode the same transformation. To reach a normal form, we start by requiring our LTWs to produce their output "as soon as possible". This method is common for transducers [2,10], and has been adapted to sequential tree-to-word transducers in [6]. In this case, the way an output word is produced by a tree-to-word can be "early" in two fashions: it can be produced sooner in the input rather than later, or it can output letters on the left of a rule rather than on the right. We take the natural extension of this definition for LTWs and find we can reuse the results and algorithms of [6].

Example 6. Consider our previous example (Ex. 2). The function $\llbracket M \rrbracket : t \mapsto 0^{|t|}$, Our transducer has only one state q , and its axiom is $ax = q$

$$\begin{aligned} q(f(x_1, x_2)) &\rightarrow 0 \cdot q(x_1) \cdot q(x_2) \\ q(a) &\rightarrow 0, \quad q(b) \rightarrow 0 \end{aligned}$$

Since all productions of q start with a 0, this LTW does not produce as up in the input as possible. To change this, we form a new state q' that produces one 0 less than q . By removing the 0 at the beginning of each rule of q , and replacing each call $q(x_i)$ by $0q'(x_i)$, we get a new equivalent LTW M' of axiom $ax' = 0 \cdot q'$

$$\begin{aligned} q'(f(x_1, x_2)) &\rightarrow 0 \cdot q'(x_1) \cdot 0 \cdot q'(x_2) \\ q'(a) &\rightarrow \varepsilon \quad q'(b) \rightarrow \varepsilon \end{aligned}$$

Example 7. Consider our previous example (Ex. 6). We could replace the first rule by $q'(f(x_1, x_2)) \rightarrow 0 \cdot 0 \cdot q'(x_1) \cdot q'(x_2)$. This new LTW would produce "more to the left", but still be equivalent to the first M .

In order to eliminate these differences in output strategies, we want transducers to produce the output as up in the input tree as possible, and then as to the left as possible. We formalize these notions in the definition of *earliest* LTWs.

Let M be a LTW, and q one of its states. We note $\text{lcp}(q)$ (respectively $\text{lcs}(q)$) the largest common prefix (respectively suffix) of the range of $\llbracket M \rrbracket_q$. We extend this definition to $\text{lcp}(qu)$, where $u \in \Delta^*$. We note $\text{lcp}(qu)$ (respectively $\text{lcs}(qu)$) the largest common prefix (respectively suffix) of the set $\{\llbracket M \rrbracket_q(t)u \mid t \in \text{dom}(\llbracket M \rrbracket_q)\}$.

Definition 8. A LTW M is earliest if it is trimmed, and:

- For every state q , $\text{lcp}(q) = \text{lcs}(q) = \varepsilon$

- For each rule $q, f \rightarrow u_0q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)})u_n \in \delta$, for every i from 1 to n , $\text{lcp}(q_i u_i) = \varepsilon$

This definition is a generalization of the one found in [6] from STWs to all LTWs. The first item ensures an earliest LTW outputs as soon as possible, the second that it produces as to the left as possible.

We will note that this second point in particular ensures a specific property on production in a rule $q, f \rightarrow u_0q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)})u_n: u_0q_1(x_{\sigma(1)}) \dots q_i(x_{\sigma(i)})u_i$ produces as much of $\llbracket M \rrbracket_q(f(s_1 \dots s_n))$ by just knowing $s_{\sigma(1)}, \dots, s_{\sigma(i)}$. We will formalize this property in the two following Lemmas. We can say that the right part of a rule cannot guess its first letter:

Lemma 9. *For M an earliest LTW, for $q, f \rightarrow u_0q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)})u_n$ one of its rules, for all i from 1 to n , $\text{lcp}(L_{q_i}u_i \dots L_{q_n}u_n) = \varepsilon$.*

Proof. This proof works by recursion, from $i = n$ to $i = 1$. For $i = n$, we know $\text{lcp}(L_{q_n}u_n) = \varepsilon$ because M is earliest.

For the recursion, for $i + 1$ we have $\text{lcp}(L_{q_{i+1}}u_{i+1} \dots L_{q_n}u_n) = \varepsilon$. Since M is earliest, we know that $\text{lcp}(q_i u_i) = \varepsilon$. We distinguish two cases:

- If there exists two letters $a \neq b$ such that $L_{q_i}u_i$ contains a word that starts with a and a word that starts with b , then $L_{q_i}u_i \dots u_n$ also contains a word that starts with a and a word that starts with b . This means $\text{lcp}(L_{q_i}u_i \dots u_n) = \varepsilon$.
- If all non-empty words of $L_{q_i}u_i$ start with the same letter a , to ensure $\text{lcp}(q_i u_i) = \varepsilon$, we have $\varepsilon \in L_{q_i}u_i$. This notably means that $L_{q_{i+1}}u_{i+1} \dots u_n \subseteq L_{q_i}u_i \dots L_{q_n}u_n$. Hence, $\text{lcp}(L_{q_i}u_i \dots L_{q_n}u_n)$ is smaller than $\text{lcp}(L_{q_{i+1}}u_{i+1} \dots L_{q_n}u_n) = \varepsilon$, which means $\text{lcp}(L_{q_i}u_i \dots L_{q_n}u_n) = \varepsilon$.

Conversely, we can say that the part on the left of the rule produces as much as possible, i.e. the lcp of all $\llbracket M \rrbracket_q(f(s_1 \dots s_n))$ for some fixed $s_{\sigma(1)}, \dots, s_{\sigma(i)}$.

Lemma 10. *For M an earliest LTW, $q, f \rightarrow u_0q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)})u_n \in \delta$, for i such that $i \leq n$, $t_{\sigma(1)}, \dots, t_{\sigma(i)}$ respectively in $\text{dom}(\llbracket M \rrbracket_{q_1}), \dots, \text{dom}(\llbracket M \rrbracket_{q_i})$, then $u_0 \llbracket M \rrbracket_{q_1}(t_{\sigma(1)}) \dots \llbracket M \rrbracket_{q_i}(t_{\sigma(i)})u_i$ is the lcp of the set:*

$$\{ \llbracket M \rrbracket_q(f(s_1, \dots, s_n)) \mid s_{\sigma(1)} = t_{\sigma(1)}, \dots, s_{\sigma(i)} = t_{\sigma(i)} \}$$

Proof. This comes from Lemma 9: for all $s = f(s_1, \dots, s_n)$, we have:

$$\llbracket M \rrbracket_q(s) = u_0 \llbracket M \rrbracket_{q_1}(t_{\sigma(1)})u_1 \dots \llbracket M \rrbracket_{q_n}(t_{\sigma(n)})u_n$$

If we fix $s_{\sigma(1)} = t_{\sigma(1)}, \dots, s_{\sigma(i)} = t_{\sigma(i)}$, all $\llbracket M \rrbracket_q(s)$ start with $u_0 \llbracket M \rrbracket_{q_1}(t_{\sigma(1)})u_1 \dots \llbracket M \rrbracket_{q_i}(t_{\sigma(i)})u_i$. This means that

$$\text{lcp}(\{ \llbracket M \rrbracket_q(f(s_1, \dots, s_n)) \mid s_{\sigma(1)} = t_{\sigma(1)}, \dots, s_{\sigma(i)} = t_{\sigma(i)} \}) = u_0 \llbracket M \rrbracket_{q_1}(t_{\sigma(1)})u_1 \dots \llbracket M \rrbracket_{q_n}(t_{\sigma(n)})u_n \text{lcp}(L_{q_{i+1}}u_{i+1} \dots u_n)$$

Since $\text{lcp}(L_{q_{i+1}}u_{i+1} \dots u_n) = \varepsilon$, we have what we wanted to prove.

Some important properties extend from [6] to earliest LTWs, most notably the fact that all LTWs can be made earliest.

Lemma 11. *For M a LTW, one can compute an equivalent earliest LTW in exponential time.*

This result is a direct generalization of the construction in Section 3 of [6]. We build the equivalent earliest LTW M' with two kinds of steps:

- If $\text{lcp}(qu) = v$, where v is a prefix of u , we can slide v through state q by creating a new state $[v^{-1}qv]$ such that for all t , $\llbracket M' \rrbracket_{[v^{-1}qv]}(t) = v^{-1}\llbracket M \rrbracket_q(t)v$. Every occurrence of $q(x_i)v$ in a rule of M is replaced by $v[v^{-1}qv](x_i)$.
 - If $\text{lcp}(q) = v$, we can produce v outside of q by creating a new state $[v^{-1}q]$ such that for all t , $\llbracket M' \rrbracket_{[v^{-1}q]}(t) = v^{-1}\llbracket M \rrbracket_q(t)$. Every occurrence of $q(x_i)$ in a rule of M is replaced by $v[v^{-1}q](x_i)$.
- Symmetrically, if $\text{lcs}(q) = v$, we create a state $[qv^{-1}]$, and every occurrence of $q(x_i)$ in a rule of M is replaced by $[qv^{-1}](x_i)v$.

Note that the exponential bound is, in fact, an exact bound, as some LTWS gain an exponential number of states through this process.

In [6], earliest STWS are actually enough to make a normal form using a Myhill-Nerode theorem: by minimizing earliest STWS (merging states with the same $\llbracket M \rrbracket_q$), we end up with a normal form with a minimal number of states. However, in the wider case of LTWS, there are still ways for two states to be equivalent and yet not syntactically equals. This impedes the process of minimization. As we will see in the next part, it remains to study how the images of subtrees can be reordered in earliest LTWS while conserving equivalence.

4 Reordering in Earliest Transducers

Syntactically different earliest LTWS may still be equivalent. Indeed, unlike sequential tree transducers [6], which impose the output to follow the order of the input, LTWS permit to flip the order.

The main point of this paper is the observation that it is sufficient to normalize the flips in the output production of earliest LTWS, in order to find a unique normal form for equivalent LTWS. To this end, we will prove that order differences are only possible in very specific cases. We start illustrating such flips in some examples, and then discuss the necessary and sufficient condition that dictates when a flip is possible.

Example 12. We reconsider Example 7. This earliest transducer “counts” the number of nodes in the input tree has only one state q' . It has the axiom $\text{ax}' = 0 \cdot q'$ and the following rules:

$$q'(f(x_1, x_2)) \rightarrow 0 \cdot 0 \cdot q'(x_1) \cdot q'(x_2), \quad q'(a) \rightarrow \varepsilon, \quad q'(b) \rightarrow \varepsilon.$$

We can now flip the order of the terms $q'(x_2)$ and $q'(x_1)$ in the first rule, and replace it by:

$$q'(f(x_1, x_2)) \rightarrow 0 \cdot 0 \cdot q'(x_2) \cdot q'(x_1).$$

This does not change $\llbracket M' \rrbracket$, since just the order is changed in which the nodes of the first and second subtree of the input are counted.

Of course, it is not always possible to flip two occurrences of terms $q_1(x_{\sigma(1)})$ and $q_2(x_{\sigma(2)})$ in LTW rules.

Example 13. Consider an earliest transducer that outputs the frontier of the input tree while replacing a by 0 and b by 1. This transducer has a single state q , the axiom $\text{ax} = q$, and the following rules:

$$q(f(x_1, x_2)) \rightarrow q(x_1) \cdot q(x_2), \quad q(a) \rightarrow 0, \quad q(b) \rightarrow 1.$$

Clearly, replacing the first rule by a flipped variant $q(f(x_1, x_2)) \rightarrow q(x_2) \cdot q(x_1)$ would not preserve transducer equivalence since $f(a, b)$ would be transformed to 10 instead of 01. More generally, no LTW with rule $q(f(x_1, x_2)) \rightarrow u_0 \cdot q_1(x_2) \cdot u_1 \cdot q_2(x_1) \cdot u_2$ produces the correct output.

Our goal is to understand the conditions when variable flips are possible.

Definition 14. For M, M' two LTWs, $q \in Q, q' \in Q'$,

$$\begin{aligned} q, f &\rightarrow u_0 q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)}) u_n \in \delta \\ q', f &\rightarrow u'_0 q'_1(x_{\sigma'(1)}) \dots q'_n(x_{\sigma'(n)}) u'_n \in \delta' \end{aligned}$$

are said to be twin rules if q and q' are equivalent.

4.1 Reordering Erasing States

We start the study of possible reordering with the obvious case of states that only produce ε : they can take every position in every rule without changing the semantics of the states. The first step towards normalization would then be to fix the positions of erasing states in the rules, to prevent differences in equivalent earliest LTWs: we put all erasing states at the end of any rule they appear in, in ascending subtree order.

Definition 15. For M a LTW, a state q is erasing if for all $t \in \text{dom}(\llbracket M \rrbracket_q)$, $\llbracket M \rrbracket_q(t) = \varepsilon$

We show that if two states are equivalent, they call erasing states on the same subtrees. We start by this length consideration:

Lemma 16. For two twin rules of earliest LTWs

$$\begin{aligned} q, f &\rightarrow u_0 q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)}) u_n \\ q', f &\rightarrow u'_0 q'_1(x_{\sigma'(1)}) \dots q'_n(x_{\sigma'(n)}) u'_n \end{aligned}$$

For i, j such that $\sigma(i) = \sigma'(j)$, and $t_{\sigma(i)} \in \text{dom}(\llbracket M \rrbracket_{q_i})$ then

$$|\llbracket M \rrbracket_{q_i}(t_{\sigma(i)})| = |\llbracket M' \rrbracket_{q'_i}(t_{\sigma(i)})|$$

Proof. The equivalence of q and q' gives for all t_1, \dots, t_n :

$$u_0 \llbracket M \rrbracket_{q_1}(t_{\sigma(1)}) \dots \llbracket M \rrbracket_{q_n}(t_{\sigma(n)}) u_n = u_0 \llbracket M' \rrbracket_{q'_1}(t_{\sigma'(1)}) \dots \llbracket M' \rrbracket_{q'_n}(t_{\sigma'(n)}) u'_n$$

We fix a value for every t_k except $t_{\sigma(i)}$. We say

$$\begin{aligned} u &= u_0 \dots \llbracket M \rrbracket_{q_{i-1}}(t_{\sigma(i-1)}) u_i, & v &= u_{i+1} \llbracket M \rrbracket_{q_{i+1}}(t_{\sigma(i+1)}) \dots u_n \\ u' &= u'_0 \dots \llbracket M' \rrbracket_{q'_{j-1}}(t_{\sigma'(j-1)}) u'_j, & v' &= u'_{j+1} \llbracket M' \rrbracket_{q'_{j+1}}(t_{\sigma'(j+1)}) \dots u'_n \end{aligned}$$

This gives us that for all $t_{\sigma(i)}$, $u\llbracket M \rrbracket_{q_i}(t_{\sigma(i)})v = u'\llbracket M' \rrbracket_{q'_j}(t_{\sigma(i)})v'$. Notably, $|\llbracket M \rrbracket_{q_i}(t_{\sigma(i)})| - |\llbracket M' \rrbracket_{q'_j}(t_{\sigma(i)})| = |u'| + |v'| - |u| - |v|$. This means the difference in size of these two production does not depend on $t_{\sigma(i)}$. We will show that this difference has to be 0.

If $|\llbracket M \rrbracket_{q_i}(t_{\sigma(i)})| > |\llbracket M' \rrbracket_{q'_j}(t_{\sigma(i)})|$ then $|u| < |u'|$, or $|v| < |v'|$. If $|u| < |u'|$, then $u' = uw$. This means $u\llbracket M \rrbracket_{q_i}(t_{\sigma(i)})v = uw\llbracket M' \rrbracket_{q'_j}(t_{\sigma(i)})v'$. For all $t_{\sigma(i)}$, $\llbracket M \rrbracket_{q_i}(t_{\sigma(i)}) \neq \varepsilon$ (it is longer than $\llbracket M' \rrbracket_{q'_j}(t_{\sigma(i)})$), and its first letter is always the first letter of w . This means $\text{lcp}(q_i) \neq \varepsilon$, which is impossible in an earliest LTW. $|v| < |v'|$ leads to $\text{lcs}(q_i) \neq \varepsilon$, another contradiction. By symmetry, $|\llbracket M' \rrbracket_{q'_j}(t_{\sigma(i)})| > |\llbracket M \rrbracket_{q_i}(t_{\sigma(i)})|$ also leads to contradiction. Therefore, both are of same size.

A direct consequence of this lemma is that in both twin rules, the subtrees read by an erasing state are the same:

Lemma 17. *For two twin rules of earliest LTWs*

$$\begin{aligned} q, f &\rightarrow u_0q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)})u_n \\ q', f &\rightarrow u'_0q'_1(x_{\sigma'(1)}) \dots q'_n(x_{\sigma'(n)})u'_n \end{aligned}$$

If q_i is erasing, and $\sigma(i) = \sigma'(j)$, then q'_j is erasing.

To normalize the order of erasing states in twin rules, we note that since erasing states produce no output letter, their position in a rule is not important to the semantics or the earliest property. We can thus push them to the right.

Lemma 18. *For M an earliest LTW, $q, f \rightarrow u_0q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)})u_n$ a rule in M , and q_i an erasing state. Then replacing this rule by*

$$q, f \rightarrow u_0q_1(x_{\sigma(1)}) \dots u_{i-1}u_i \dots q_n(x_{\sigma(n)})u_nq_i(x_{\sigma(i)})$$

does not change $\llbracket M \rrbracket_q$, and M remains earliest.

Proof. While the proof that $\llbracket M \rrbracket_q$ remains unchanged is technically an inductive proof to show that M and all its states keep the same semantics, it is clear that the only important step is to ensure that for $s = f(s_1, \dots, s_n)$, $\llbracket M \rrbracket_q(s)$ remains unchanged. With the old rule, we have $\llbracket M \rrbracket_q(s) = u_0 \dots u_{i-1} \llbracket M \rrbracket_{q_i}(s_{\sigma(i)}) u_i \dots u_n$. Since q_i is erasing, $\llbracket M \rrbracket_{q_i}(s_{\sigma(i)}) = \varepsilon$ and $\llbracket M \rrbracket_q(s) = u_0 \dots u_{i-1} u_i \dots u_n \llbracket M \rrbracket_{q_i}(s_{\sigma(i)})$. Hence, the rule $q, f \rightarrow u_0q_1(x_{\sigma(1)}) \dots u_{i-1}u_i \dots q_n(x_{\sigma(n)})u_nq_i(x_{\sigma(i)})$ produces the right output: $\llbracket M \rrbracket_q$ does not change.

To prove that M remains earliest, we just have to check that $\text{lcp}(q_{i-1}u_{i-1}u_i) = \varepsilon$. Since M was earliest with the previous rule, we have $\text{lcp}(q_iu_i) = \varepsilon$. Since q_i is erasing, this means $u_i = \varepsilon$. Since M was earliest with the previous rule, we have $\text{lcp}(q_{i-1}u_{i-1}) = \varepsilon$. Hence, $\text{lcp}(q_{i-1}u_{i-1}u_i) = \varepsilon$.

Given this lemma, we can define a first normalization step where all erasing states appear at the end of the rules in ascending subtree order.

Definition 19. *An earliest LTW M is erase-ordered if for every rule*

$q, f \rightarrow u_0q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)})u_n \in \delta$, if q_i is erasing, then for all $j > i$, q_j is erasing, and $\sigma(i) < \sigma(j)$.

Lemma 20. *For M an earliest LTW, one can make M erase-ordered in polynomial time without changing the semantic of its states.*

Proof. We can detect which states are erasing in polynomial time: as shown in Lemma 3, we can compute in polynomial time a grammar that describes L_q . Deciding if a grammar describes $\{\varepsilon\}$ can be decided in polynomial time. From there, Lemma 18 ensures that making a LTW erase-ordered is just a matter of pushing all erasing states at the end of the rules and then sorting them in ascending subtree order, which can be done in polynomial time.

4.2 Reordering Producing States

As we saw in Example 13, some flips between states are not possible. We will now study what makes reordering non-erasing states possible. As we will see, only few differences are possible between twin rules in erase-ordered earliest LTWs. Two states transforming the same subtree are equivalent, and the only order differences are caused by flipping states whose productions commute in Δ^* .

To prove this, we begin by establishing a few preliminary results. We first show that to the left of σ and σ' 's first difference, both rules are identical.

Lemma 21. *For two twin rules of erase-ordered earliest LTWs M, M'*

$$\begin{aligned} q, f &\rightarrow u_0 q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)}) u_n \\ q', f &\rightarrow u'_0 q'_1(x_{\sigma'(1)}) \dots q'_n(x_{\sigma'(n)}) u'_n \end{aligned}$$

For i such that if $k \leq i$ then $\sigma(k) = \sigma'(k)$, $\llbracket M \rrbracket_{q_i} = \llbracket M' \rrbracket_{q'_i}$, and $u_i = u'_i$.

Proof. This results from Lemma 10: if σ and σ' coincide before i , then for all $t_{\sigma(1)}, \dots, t_{\sigma(i)}$, $u_0 \llbracket M \rrbracket_{q_1}(t_{\sigma(1)}) \dots u_i$ and $u'_0 \llbracket M' \rrbracket_{q'_1}(t_{\sigma'(1)}) \dots u'_i$ are both equal to the lcp of $\{\llbracket M \rrbracket_q(f(s_1, \dots, s_n)) \mid s_{\sigma(1)} = t_{\sigma(1)}, \dots, s_{\sigma(n)} = t_{\sigma(n)}\}$. This means that:

$$u_0 \llbracket M \rrbracket_{q_1}(t_{\sigma(1)}) \dots \llbracket M \rrbracket_{q_i}(t_{\sigma(i)}) u_i = u'_0 \llbracket M' \rrbracket_{q'_1}(t_{\sigma'(1)}) \dots \llbracket M' \rrbracket_{q'_i}(t_{\sigma'(i)}) u'_i$$

Since this is also true for $i-1$, we can remove everything but the last part for each side of this equation, to obtain that for all $t_{\sigma(i)}$, $\llbracket M \rrbracket_{q_i}(t_{\sigma(i)}) u_i = \llbracket M' \rrbracket_{q'_i}(t_{\sigma(i)}) u'_i$. Lemma 16 gives us $|\llbracket M \rrbracket_{q_i}(t_{\sigma(i)})| = |\llbracket M' \rrbracket_{q'_i}(t_{\sigma'(i)})|$, and $u_i = u'_i$. This means that q_i and q'_i are equivalent, and $u_i = u'_i$.

It still remains to show what happens when σ and σ' stop coinciding. We study the leftmost order difference between two twin rules in erasing-ordered earliest LTWs, that is to say the smallest i such that $\sigma(i) \neq \sigma'(i)$. Note that Lemma 17 ensures that such a difference occurs before the end of the rule where the erasing states are sorted.

Lemma 22. *For two twin rules of erase-ordered earliest LTWs M, M'*

$$\begin{aligned} q, f &\rightarrow u_0 q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)}) u_n \\ q', f &\rightarrow u'_0 q'_1(x_{\sigma'(1)}) \dots q'_n(x_{\sigma'(n)}) u'_n \end{aligned}$$

For i such that $\sigma(i) \neq \sigma'(i)$ and for any $k < i$, $\sigma(k) = \sigma'(k)$, for j such that $\sigma'(i) = \sigma(j)$, we have:

- (A) For all k from i to $j-1$, $u_k = \varepsilon$ and there exists $t_{\sigma(k)}^\varepsilon$
such that $\llbracket M \rrbracket_{q_k}(t_{\sigma(k)}^\varepsilon) = \varepsilon$
- (A') For all k from i to j , for k' such that $\sigma(k) = \sigma'(k')$, for l from i to $k'-1$,
 $u'_l = \varepsilon$ and there exists $t_{\sigma'(l)}^\varepsilon$ such that $\llbracket M' \rrbracket_{q'_l}(t_{\sigma'(l)}^\varepsilon) = \varepsilon$
- (B) For all k from i to j , for k' such that $\sigma(k) = \sigma'(k')$, q_k is equivalent to $q'_{k'}$
- (C) All q_i, \dots, q_j are periodic of same period.

Proof. We first prove point (A), then (A') with the same arguments. We then use them to show point (B), then from (A), (A') and (B) we finally show point (C).

For point (A), we use the equivalence of q and q' . For all t_1, \dots, t_n ,

$$u_0 \llbracket M \rrbracket_{q_1}(t_{\sigma(1)}) \dots \llbracket M \rrbracket_{q_n}(t_{\sigma(n)}) u_n = u_0 \llbracket M' \rrbracket_{q'_1}(t_{\sigma'(1)}) \dots \llbracket M' \rrbracket_{q'_n}(t_{\sigma'(n)}) u'_n$$

Lemma 21 gives us that everything up to u_{i-1} and u'_{i-1} coincide. We then get

$$\llbracket M \rrbracket_{q_i}(t_{\sigma(i)}) \dots \llbracket M \rrbracket_{q_n}(t_{\sigma(n)}) u_n = \llbracket M' \rrbracket_{q'_i}(t_{\sigma'(i)}) \dots \llbracket M' \rrbracket_{q'_n}(t_{\sigma'(n)}) u'_n$$

Since q'_i is not erasing, we can fix $t_{\sigma'(i)}$ such that $\llbracket M' \rrbracket_{q'_i}(t_{\sigma'(i)}) \neq \varepsilon$. We call its first letter a . All non- ε productions of q_i must begin by a . This is only possible in an earliest if there exists $t_{\sigma(i)}^\varepsilon$ such that $\llbracket M \rrbracket_{q_i}(t_{\sigma(i)}^\varepsilon) = \varepsilon$. We now fix $t_{\sigma(i)} = t_{\sigma(i)}^\varepsilon$. If $u_i \neq \varepsilon$, its first letter is a . This is impossible in an earliest since it would mean $\text{lcp}(q_i u_i) \neq \varepsilon$. Hence $u_i = \varepsilon$. We can make the same reasoning for q_{i+1} and u_{i+1} , and so on all the way to q_{j-1} and u_{j-1} .

Point (A') uses a similar argument. we start from

$$\llbracket M \rrbracket_{q_i}(t_{\sigma(i)}) \dots \llbracket M \rrbracket_{q_n}(t_{\sigma(n)}) u_n = \llbracket M' \rrbracket_{q'_i}(t_{\sigma'(i)}) \dots \llbracket M' \rrbracket_{q'_n}(t_{\sigma'(n)}) u'_n$$

Using point (A), we choose $t_{\sigma(i)} = t_{\sigma(i)}^\varepsilon, \dots, t_{\sigma(k-1)} = t_{\sigma(k-1)}^\varepsilon$. We get

$$\llbracket M \rrbracket_{q_k}(t_{\sigma(k)}) \dots \llbracket M \rrbracket_{q_n}(t_{\sigma(n)}) u_n = \llbracket M' \rrbracket_{q'_i}(t_{\sigma'(i)}) \dots \llbracket M' \rrbracket_{q'_n}(t_{\sigma'(n)}) u'_n$$

Using the same argument as we used for point (A), we can choose $t_{\sigma(k)}$ such that $\llbracket M \rrbracket_{q_k}(t_{\sigma(k)}) \neq \varepsilon$, and conclude that for l from i to $k'-1$, $u'_l = \varepsilon$ and there exists $t_{\sigma'(l)}^\varepsilon$ such that $\llbracket M' \rrbracket_{q'_l}(t_{\sigma'(l)}^\varepsilon) = \varepsilon$.

For point (B), we use point (A) and (A') to eliminate everything in front of q_k and $q'_{k'}$ by picking all $t_{\sigma(l)}^\varepsilon$ up to $k-1$ and all $t_{\sigma'(l')}^\varepsilon$ up to $k'-1$.

$$\llbracket M \rrbracket_{q_k}(t_{\sigma(k)}) \dots \llbracket M \rrbracket_{q_n}(t_{\sigma(n)}) u_n = \llbracket M' \rrbracket_{q'_{k'}}(t_{\sigma'(k')}) \dots \llbracket M' \rrbracket_{q'_n}(t_{\sigma'(n)}) u'_n$$

From Lemma 16, we know that $|\llbracket M \rrbracket_{q_k}(t_{\sigma(k)})| = |\llbracket M' \rrbracket_{q'_{k'}}(t_{\sigma'(k')})|$. We conclude that q_k and $q'_{k'}$ are equivalent.

For point (C), we take k' such that $\sigma(k) = \sigma'(k')$. We use (A) to erase everything but q_k, q_j, q'_i and $q'_{k'}$ by picking every $t_{\sigma(l)}^\varepsilon$ and $t_{\sigma'(l')}^\varepsilon$ except theirs.

$$\llbracket M \rrbracket_{q_k}(t_{\sigma(k)}) \llbracket M \rrbracket_{q_j}(t_{\sigma(j)}) \dots u_n = \llbracket M' \rrbracket_{q'_i}(t_{\sigma'(i)}) \llbracket M' \rrbracket_{q'_{k'}}(t_{\sigma'(k')}) \dots u'_n$$

Point (B) gives q_k is equivalent to $q'_{k'}$ and q_j is equivalent to q'_i . We get that $\llbracket M \rrbracket_{q_k}(t_{\sigma(k)}) \llbracket M \rrbracket_{q_j}(t_{\sigma(j)}) = \llbracket M \rrbracket_{q_j}(t_{\sigma(j)}) \llbracket M \rrbracket_{q_k}(t_{\sigma(k)})$. This means that the productions of q_k and q_j commute, which in Δ^* is equivalent to say they are words of same period. Therefore, q_j and q_k are periodic of same period.

This result allows us to resolve the first order difference between two twin rules by flipping q_j with neighbouring periodic states of same period. We can iterate this method to solve all order differences.

Theorem 23. For two twin rules of erase-ordered earliest LTWs,

$$\begin{aligned} q, f &\rightarrow u_0 q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)}) u_n \\ q', f &\rightarrow u'_0 q'_1(x_{\sigma'(1)}) \dots q'_n(x_{\sigma'(n)}) u'_n \end{aligned}$$

One can replace the rule of q to another rule of same subtree order as the rule of q' only by flipping neighbour states q_k and q_{k+1} of same period where $u_k = \varepsilon$.

We can use Lemma 22 to solve the leftmost difference: for i first index such that $\sigma(i) \neq \sigma'(i)$, and j such that $\sigma(i) = \sigma'(j)$, we have $u_i = \dots = u_{j-1} = \varepsilon$ and q_i, \dots, q_j commute with each other. This means we can replace the first rule by:

$$q, f \rightarrow u_0 \dots q_j(x_{\sigma(j)}) q_i(x_{\sigma(i)}) \dots q_{j-1}(x_{\sigma(j-1)}) u_j \dots u_n$$

where $q_j(x_{\sigma(j)})$ is to the left of $q_i(x_{\sigma(i)}) \dots q_{j-1}(x_{\sigma(j-1)})$ without changing $\llbracket M \rrbracket_q$.

This solves the leftmost order difference: we can iterate this method until both rules have the same order.

Finally, we call Lemma 21 on the rules reordered by Theorem 23 to show that two twin rules use equivalent states and the same constant words:

Theorem 24. For two twin rules of erase-ordered earliest LTWs,

$$\begin{aligned} q, f &\rightarrow u_0 q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)}) u_n \\ q', f &\rightarrow u'_0 q'_1(x_{\sigma'(1)}) \dots q'_n(x_{\sigma'(n)}) u'_n \end{aligned}$$

$u_0 = u'_0, \dots, u_n = u'_n$, and for k, k' such that $\sigma(k) = \sigma'(k')$, $\llbracket M \rrbracket_{q_k} = \llbracket M' \rrbracket_{q'_{k'}}$.

5 Myhill-Nerode Theorem and Normal Form

In Section 3, we showed that LTWs can be made earliest. In Section 4, we first showed that all earliest LTWs can be made erase-ordered, then we made explicit what reorderings are possible between two rules of two equivalent states. In this section, we use these results to fix a reordering strategy. This will give us a new normal form, *ordered earliest* LTWs. We will show that each LTW in equivalent to a unique minimal ordered earliest LTW, whose size is at worst exponential.

We first use Theorem 23 to define a new normal form: ordered earliest LTWs.

Definition 25. A LTW M is said to be ordered earliest if it is earliest, and for each rule $q, f \rightarrow u_0 q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)}) u_n$:

- If q_i is erasing, then for any $j > i$, q_j is erasing.
- If $u_i = \varepsilon$, and q_i and q_{i+1} are periodic of same period, $\sigma(i) < \sigma(i+1)$.

Note that this definition notably implies that any ordered earliest is erase-ordered earliest. On top of that, we impose that if two adjacent states are periodic of same period, and thus could be flipped, they are sorted by ascending subtree.

Lemma 26. For M an earliest LTW, one can make M ordered in polynomial time without changing the semantic of its states.

Proof. We saw in Lemma 20 that one can push and sort erasing states. For this result, sorting periodic states is not more complicated. However, one must test first whether two states are periodic of same period. This can be done in polynomial time. Lemma 3 shows how to build a CFG whose language is the production of a state. Then, the problem of deciding if two CFG languages are periodic of same period is known to be polynomial (see for example [12,9]).

These ordered LTWs have a very important property for normalization purposes: two twin rules have the same order σ .

Lemma 27. *For two twin rules of ordered earliest LTWs,*

$$\begin{aligned} q, f &\rightarrow u_0 q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)}) u_n \\ q', f &\rightarrow u'_0 q'_1(x_{\sigma'(1)}) \dots q'_n(x_{\sigma'(n)}) u'_n \end{aligned}$$

$u_0 = u'_0, \dots, u_n = u'_n, \sigma = \sigma'$ and for all k , $\llbracket M \rrbracket_{q_k} = \llbracket M' \rrbracket_{q'_k}$.

Proof. Theorem 24 gives everything if we can prove $\sigma = \sigma'$. To prove $\sigma = \sigma'$, we suppose there is a difference, then find a contradiction with the fact that both transducers are ordered. Suppose $\sigma \neq \sigma'$. We call i the smallest index such that $\sigma(i) \neq \sigma'(i)$. Lemma 22 tells us that for j such that $\sigma(j) = \sigma'(i)$, all q_i, q_{i+1}, \dots, q_j are periodic of same period. By symmetry, for j' such that $\sigma(i) = \sigma'(j')$, all $q'_i, q'_{i+1}, \dots, q'_{j'}$ are periodic of same period. If $\sigma(i) > \sigma'(i)$, then $\sigma(i) > \sigma(j)$ and the rule of q cannot belong to an ordered LTW, since q_j should appear before q_i . If $\sigma'(i) > \sigma(i)$, then $\sigma'(i) > \sigma'(j')$ and the rule of q' cannot belong to an ordered LTW, since $q'_{j'}$ should appear before q'_i . This is a contradiction: we conclude that it is impossible for σ to be different to σ' .

Our goal is now to show the existence of a unique minimal normal LTW equivalent to any M . We start by showing that if we run in parallel two equivalent earliest LTWs, they use equivalent states when transforming the same subtree. We first formalize this intuition by defining co-reachable states, and their desired property:

Definition 28. *For two LTWs M, M' , we define co-reachable pair of $Q \times Q'$:*

- If $ax = u_0 q_0 v_0$ and $ax' = u'_0 q'_0 v'_0$ then q_0 and q'_0 are co-reachable.
- If q and q' are co-reachable,

$$\begin{aligned} q, f &\rightarrow u_0 q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)}) u_n \\ q', f &\rightarrow u'_0 q'_1(x_{\sigma'(1)}) \dots q'_n(x_{\sigma'(n)}) u'_n \end{aligned}$$

and k, k' such that $\sigma(k) = \sigma'(k')$, then q_k and $q'_{k'}$ are co-reachable.

Lemma 29. *For two equivalent earliest LTWs M and M' , if q and q' are co-reachable, then $\llbracket M \rrbracket_q = \llbracket M' \rrbracket_{q'}$*

Proof. By induction on co-reachable pairs. If $ax = u_0 q_0 v_0$ and $ax' = u'_0 q'_0 v'_0$, since M and M' are earliest, $u_0 = \text{lcp}(\llbracket M \rrbracket) = \text{lcp}(\llbracket M' \rrbracket) = u'_0$. Then, $v_0 = \text{lcs}(q_0 v_0) = \text{lcs}(q'_0 v'_0) = v'_0$. We then get that q_0 and q'_0 are equivalent.

The inductive case uses Theorem 24: for q_i and q'_j co-reachable such that there exists q, q' co-reachable (and equivalent by induction hypothesis), and two twin rules

$$\begin{aligned} q, f &\rightarrow u_0 q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)}) u_n \\ q', f &\rightarrow u'_0 q'_1(x_{\sigma'(1)}) \dots q'_n(x_{\sigma'(n)}) u'_n \end{aligned}$$

such that $\sigma(i) = \sigma(j)$, Theorem 24 ensures that q_i and q'_j are equivalent.

This allows us to show that two earliest LTWS use equivalent states:

Lemma 30. *For two equivalent earliest LTWS M and M' , for q state of M , there exist an equivalent state q' in M'*

Proof. Since M is earliest, it is trimmed. This means that any state q of M can be reached in a run of M . Therefore, by making the same run in parallel in M' , we find q' such that q and q' are co-reachable, and by Lemma 29, equivalent.

Since all equivalent earliest LTWS use the same states, they have the *minimal* amount of states when they don't have two redundant states q, q' such that $\llbracket M \rrbracket_q = \llbracket M \rrbracket_{q'}$.

We first show that such a minimal ordered LTW can always be found from an ordered LTW efficiently:

Lemma 31. *For M an ordered LTW, one can make an equivalent M' minimal ordered in polynomial time*

Proof. The idea behind the algorithm is simple: if two states q and q' are equivalent, they are "merged": q' is deleted, and every occurrence of q' in rules is replaced by an occurrence of q . It is clear that such a merge leaves M' equivalent to M . Furthermore, when no new redundant states are found, M' is minimal.

This algorithm can only work if we can test in polynomial time if two states are redundant. To this end, we note that Lemma 27 gives a necessary and sufficient condition for two states to be equivalent: $\llbracket M \rrbracket_q = \llbracket M \rrbracket_{q'}$ if and only if for every rule $q, f \rightarrow u_0 q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)}) u_n$ there exists a matching rule $q', f \rightarrow u'_0 q'_1(x_{\sigma'(1)}) \dots q'_n(x_{\sigma'(n)}) u'_n$ such that $u_1 = u'_1, \dots, u_n = u'_n, \sigma = \sigma'$, for all k from 1 to n , $\llbracket M \rrbracket_{q_k} = \llbracket M \rrbracket_{q'_k}$, and reciprocally every rule of q' has a matching rule in q .

From there finding out which states are redundant can be done using an algorithm similar to the minimization of deterministic finite word automata: we find the equivalent classes of indistinguishable states by supposing them all indistinguishable at first, then refining this assumption through examination of their rules using the previously described criterion. When the refining ends, only indistinguishable (and therefore equivalent) states remain in the same equivalence class.

We remind that this means this minimal normal form can be found from any LTW in exponential time, as a combination of Lemma 11, Lemma 26, and Lemma 31:

Theorem 32. *For M a LTW, one can make an equivalent M' minimal ordered in exponential time*

All that remains is to show that minimal ordered LTWs characterise a *unique minimal normal form* on LTWs.

Theorem 33. *For M a LTW, there exists a unique minimal ordered earliest LTW M' equivalent to M (up to state renaming).*

Proof. The existence of such a minimal ordered earliest LTW derives directly from Theorem 32.

The uniqueness derives from several properties we showed in this paper. Imagine M and M' two equivalent minimal ordered earliest LTWs. The fact that they have equivalent states come from Lemma 30. Since both are minimal, neither have redundant state: each q of M is equivalent to exactly one q' of M' and vice-versa. We will now show that any rule of M' is a rule of M where every occurrence of a state of M is replaced by its unique equivalent match in M' . Consider two twin rules in these transducers.

$$\begin{aligned} q, f &\rightarrow u_0 q_1(x_{\sigma(1)}) \dots q_n(x_{\sigma(n)}) u_n \\ q', f &\rightarrow u'_0 q'_1(x_{\sigma'(1)}) \dots q'_n(x_{\sigma'(n)}) u'_n \end{aligned}$$

Lemma 27 ensures that $u_0 = u'_0, \dots, u_n = u'_n, \sigma = \sigma'$ and for all k , $\llbracket M \rrbracket_{q_k} = \llbracket M' \rrbracket_{q'_k}$. Since M and M' are minimal, each occurrence of q_k is replaced by the same q'_k , the unique state of M' equivalent to q_k . In other words, both rules are identical up to renaming each state of M by its unique equivalent state of M' , or vice-versa.

6 Conclusion and Future Work

This paper's goal was to solve the equivalence problem on linear tree-to-word transducers, by establishing a normal form and a Myhill-Nerode theorem on this class. To do so we naturally extended the notion of earliest transducers that already existed in sequential tree transducers [6]. However it appeared that this was no longer enough to define a normal form: we studied all possible reorderings that could happen in an earliest LTW. We then used this knowledge to define a new normal form, that has both an output strategy (earliest) and an ordering strategy (ordered earliest), computable from any LTW in EXPTIME.

There are several ways to follow up on this result: one would be adapting the learning algorithm presented in [7], accounting for the fact that we now also have to learn the order in which the images appear. It could also be relevant to note that in [6], another algorithm decides equivalence in polynomial time, which is more efficient than computing the normal form. Such an algorithm would be an improvement over the actual randomized polynomial algorithm by [13]. As far as Myhill-Nerode theorems go, the next step would be to consider all tree-to-word transducers. This problem is known to be difficult. Recently, [13] gave a randomized polynomial algorithm to decide equivalence, but did not provide a Myhill-Nerode characterization.

References

1. Rajeev Alur and Loris D'Antoni. Streaming tree transducers. *CoRR*, abs/1104.2599, 2011.
2. Christian Choffrut. Minimizing subsequential transducers: a survey. *Theor. Comput. Sci.*, 292(1):131–143, 2003.
3. Joost Engelfriet and Sebastian Maneth. Macro tree translations of linear size increase are MSO definable. *SIAM J. Comput.*, 32(4):950–1006, 2003.
4. Joost Engelfriet and Sebastian Maneth. The equivalence problem for deterministic MSO tree transducers is decidable. In *FSTTCS 2005: Foundations of Software Technology and Theoretical Computer Science, 25th International Conference, Hyderabad, India, December 15-18, 2005, Proceedings*, pages 495–504, 2005.
5. Joost Engelfriet and Heiko Vogler. Macro tree transducers. *J. Comput. Syst. Sci.*, 31(1):71–146, 1985.
6. Grégoire Laurence, Aurélien Lemay, Joachim Niehren, Slawek Staworko, and Marc Tommasi. Normalization of sequential top-down tree-to-word transducers. In *Language and Automata Theory and Applications - 5th International Conference, LATA 2011, Tarragona, Spain, May 26-31, 2011. Proceedings*, pages 354–365, 2011.
7. Grégoire Laurence, Aurélien Lemay, Joachim Niehren, Slawek Staworko, and Marc Tommasi. Learning sequential tree-to-word transducers. In *Language and Automata Theory and Applications - 8th International Conference, LATA 2014, Madrid, Spain, March 10-14, 2014. Proceedings*, pages 490–502, 2014.
8. Aurélien Lemay, Sebastian Maneth, and Joachim Niehren. A learning algorithm for top-down XML transformations. In *Proceedings of the Twenty-Ninth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS 2010, June 6-11, 2010, Indianapolis, Indiana, USA*, pages 285–296, 2010.
9. Markus Lohrey. *The Compressed Word Problem for Groups*. Springer, 2014.
10. Sebastian Maneth and Helmut Seidl. Deciding equivalence of top-down XML transformations in polynomial time. In *PLAN-X 2007, Programming Language Technologies for XML, An ACM SIGPLAN Workshop colocated with POPL 2007, Nice, France, January 20, 2007*, pages 73–79, 2007.
11. José Oncina, Pedro Garcia, and Enrique Vidal. Learning subsequential transducers for pattern recognition interpretation tasks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 15(5):448–458, 1993.
12. Wojciech Plandowski. The complexity of the morphism equivalence problem for context-free languages. *Ph. D. thesis, Dept. of Mathematics, Informatics and Mechanics*, 1995.
13. Helmut Seidl, Sebastian Maneth, and Gregor Kemper. Equivalence of deterministic top-down tree-to-string transducers is decidable. *CoRR*, abs/1503.09163, 2015.