



# The Bottom-Up Position Tree Automaton and Its Compact Version

Samira Attou<sup>1</sup>, Ludovic Mignot<sup>2(✉)</sup>, and Djelloul Ziadi<sup>2</sup>

<sup>1</sup> Faculty of Mathematics, RECITS Laboratory, USTHB,  
BP 32, El Alia, 16111 Bab Ezzouar, Algiers, Algeria  
sattou@usthb.dz

<sup>2</sup> Groupe de Recherche Rouennais en Informatique Fondamentale,  
Université de Rouen Normandie,  
Avenue de l'Université, 76801 Saint-Étienne-du-Rouvray, France  
{ludovic.mignot,djelloul.ziadi}@univ-rouen.fr

**Abstract.** The conversion of a given regular tree expression into a tree automaton has been widely studied. However, classical interpretations are based upon a Top-Down interpretation of tree automata. In this paper, we propose a new construction based on Glushkov's one using a Bottom-Up interpretation. One of the main goals of this technique is to consider as a next step the links with deterministic recognizers, consideration that cannot be performed with classical Top-Down approaches. Furthermore, we exhibit a method to factorize transitions of tree automata and show that this technique is particularly interesting for the Glushkov constructions, by considering natural factorizations due to the structure of regular expression.

## 1 Introduction

Automata are recognizers used in various domains of applications especially in computer science, *e.g.* to represent (non necessarily finite) languages, or to solve the membership test, *i.e.* to verify whether a given element belongs to a language or not. Regular expressions are compact representations for these recognizers. Indeed, in the case where elements are words, it is well known that each regular expression can be transformed into a finite state machine recognizing the language it defines. Several methods have been proposed to realize this conversion. As an example, Glushkov [6] (and independently Mc-Naughton and Yamada [9]) showed how to construct a non deterministic finite automaton with  $n + 1$  states where  $n$  represents the number of letters of a given regular expression. The main idea of the construction is to define some particular sets named First, Follow and Last that are computed with respect to the occurrences of the symbols that appear in the expression.

These so-called Glushkov automata (or position automata) are finite state machines that have been deeply studied. They have been structurally characterized by Caron and Ziadi [4], allowing us to invert the Glushkov computation

by constructing an expression with  $n$  symbols from a Glushkov automaton with  $n + 1$  states. They have been considered too in the theoretical notion of one-unambiguity by Bruggemann-Klein and Wood [3], characterizing regular languages recognized by a deterministic Glushkov automaton, or with practical thoughts, like expression updating [2]. Finally, it is also related to combinatorial research topics. As an example, Nicaud [12] proved that the average number of transitions of Glushkov automata is linear.

The Glushkov construction was extended to tree automata [8,11], using a Top-Down interpretation of tree expressions. This interpretation can be problematic while considering determinism. Indeed, it is a folklore that there exist regular tree languages that cannot be recognized by Top-Down deterministic tree automata. Extensions of one-ambiguity are therefore incompatible with this approach.

In this paper, we propose a new approach based on the construction of Glushkov in a Bottom-Up interpretation. We also define a compressed version of tree automata in order to factorize the transitions, and we show how to apply it directly over the Glushkov computation using natural factorizations due to the structure of the expressions. The paper is structured as follows: in Sect. 2, we recall some properties related to regular tree expressions; we also introduce some basic definitions. We define, in Sect. 3, the position functions used for the construction of the Bottom-Up position tree automaton. Section 4 indicates the way that we construct the Bottom-Up position tree automaton with a linear number of states using the functions shown in Sect. 3. In Sect. 5, we propose the notion of compressed automaton and show how to reduce the size of the position automaton computed in the previous section.

## 2 Preliminaries

Let us first introduce some notations and preliminary definitions. For a boolean condition  $\psi$ , we denote by  $(E \mid \psi)$   $E$  if  $\psi$  is satisfied,  $\emptyset$  otherwise. Let  $\Sigma = (\Sigma_n)_{n \geq 0}$  be a finite ranked alphabet. A *tree*  $t$  over  $\Sigma$  is inductively defined by  $t = f(t_1, \dots, t_k)$  where  $f \in \Sigma_k$  and  $t_1, \dots, t_k$  are  $k$  trees over  $\Sigma$ . The relation “ $s$  is a subtree of  $t$ ” is denoted by  $s \prec t$  for any two trees  $s$  and  $t$ . We denote by  $\text{root}(t)$  the root symbol of the tree  $t$ , *i.e.*

$$\text{root}(f(t_1, \dots, t_k)) = f. \quad (1)$$

The *predecessors* of a symbol  $f$  in a tree  $t$  are the symbols that appear directly above it. We denote by  $\text{father}(t, f)$ , for a tree  $t$  and a symbol  $f$  the pairs

$$\text{father}(t, f) = \{(g, i) \in \Sigma_l \times \mathbb{N} \mid \exists g(s_1, \dots, s_l) \prec t, \text{root}(s_i) = f\}. \quad (2)$$

These couples link the predecessors of  $f$  and the indices of the subtrees in  $t$  that  $f$  is the root of. Let us consider a tree  $t = g(t_1, \dots, t_k)$  and a symbol  $f$ . By definition of the structure of a tree, a predecessor of  $f$  in  $t$  is a predecessor of  $f$

in a subtree  $t_i$  of  $t$ , or  $g$  if  $f$  is a root of a subtree  $t_i$  of  $t$ . Consequently:

$$\text{father}(t, f) = \bigcup_{i \leq n} \text{father}(t_i, f) \cup \{(g, i) \mid f \in \text{root}(t_i)\}. \quad (3)$$

We denote by  $T_\Sigma$  the set of trees over  $\Sigma$ . A tree language  $L$  is a subset of  $T_\Sigma$ .

For any 0-ary symbol  $c$ , let  $t \cdot_c L$  denote the tree language constituted of the trees obtained by substitution of any symbol  $c$  of  $t$  by a tree of  $L$ . By a linear extension, we denote by  $L \cdot_c L' = \{t \cdot_c L' \mid t \in L\}$ . For an integer  $n$ , the  $n$ -th substitution  ${}^{c,n}$  of a language  $L$  is the language  $L^{c,n}$  recursively defined by

$$L^{c,n} = \begin{cases} \{c\}, & \text{if } n = 0, \\ L \cdot_c L^{c,n-1} & \text{otherwise.} \end{cases}$$

Finally, we denote by  $L(E_1^{*c})$  the language  $\bigcup_{k \geq 0} L(E_1)^{c,k}$ .

An *automaton* over  $\Sigma$  is a 4-tuple  $A = (\bar{Q}, \Sigma, Q_F, \delta)$  where  $Q$  is a set of states,  $Q_F \subseteq Q$  is the set of final states, and  $\delta \subseteq \bigcup_{k \geq 0} (Q^k \times \Sigma_k \times Q)$  is the set of transitions, which can be seen as the function from  $Q^k \times \Sigma_k$  to  $2^Q$  defined by

$$(q_1, \dots, q_k, f, q) \in \delta \Leftrightarrow q \in \delta(q_1, \dots, q_k, f).$$

It can be linearly extended as the function from  $(2^Q)^k \times \Sigma_k$  to  $2^Q$  defined by

$$\delta(Q_1, \dots, Q_n, f) = \bigcup_{(q_1, \dots, q_n) \in Q_1 \times \dots \times Q_n} \delta(q_1, \dots, q_n, f). \quad (4)$$

Finally, we also consider the function  $\Delta$  from  $T_\Sigma$  to  $2^Q$  defined by

$$\Delta(f(t_1, \dots, t_n)) = \delta(\Delta(t_1), \dots, \Delta(t_n), f).$$

Using these definitions, the language recognized by the automaton  $A$  is the language  $\{t \in T_\Sigma \mid \Delta(t) \cap Q_F \neq \emptyset\}$ .

A *regular expression*  $E$  over the alphabet  $\Sigma$  is inductively defined by:

$$\begin{aligned} E &= f(E_1, \dots, E_k), & E &= E_1 + E_2, \\ E &= E_1 \cdot_c E_2, & E &= E_1^{*c}, \end{aligned}$$

where  $k \in \mathbb{N}$ ,  $c \in \Sigma_0$ ,  $f \in \Sigma_k$  and  $E_1, \dots, E_k$  are any  $k$  regular expressions over  $\Sigma$ . In what follows, we consider expressions where the subexpression  $E_1 \cdot_c E_2$  only appears when  $c$  appears in the expression  $E_1$ . The *language denoted* by  $E$  is the language  $L(E)$  inductively defined by

$$\begin{aligned} L(f(E_1, \dots, E_k)) &= \{f(t_1, \dots, t_k) \mid t_j \in L(E_j), j \leq k\}, \\ L(E_1 + E_2) &= L(E_1) \cup L(E_2), \\ L(E_1 \cdot_c E_2) &= L(E_1) \cdot_c L(E_2), \\ L(E_1^{*c}) &= L(E_1)^{*c}, \end{aligned}$$

with  $k \in \mathbb{N}$ ,  $c \in \Sigma_0$ ,  $f \in \Sigma_k$  and  $E_1, \dots, E_k$  any  $k$  regular expressions over  $\Sigma$ .

A regular expression  $E$  is *linear* if each symbol  $\Sigma_n$  with  $n \neq 0$  occurs at most once in  $E$ . Note that the symbols of rank 0 may appear more than once. We denote by  $\bar{E}$  the linearized form of  $E$ , which is the expression  $E$  where any occurrence of a symbol is indexed by its position in the expression. The set of indexed symbols, called *positions*, is denoted by  $\text{Pos}(\bar{E})$ . We also consider the *delinearization* mapping  $h$  sending a linearized expression over its original unindexed version.

Let  $\phi$  be a function between two alphabets  $\Sigma$  and  $\Sigma'$  such that  $\phi$  sends  $\Sigma_n$  to  $\Sigma'_n$  for any integer  $n$ . By a well-known adjunction, this function is extended to an *alphabetical morphism* from  $T(\Sigma)$  to  $T(\Sigma')$  by setting  $\phi(f(t_1, \dots, t_n)) = \phi(f)(\phi(t_1), \dots, \phi(t_n))$ . As an example, one can consider the delinearization morphism  $h$  that sends an indexed alphabet to its unindexed version. Given a language  $L$ , we denote by  $\phi(L)$  the set  $\{\phi(t) \mid t \in L\}$ . The *image by  $\phi$*  of an automaton  $A = (\Sigma, Q, Q_F, \delta)$  is the automaton  $\phi(A) = (\Sigma', Q, Q_F, \delta')$  where

$$\delta' = \{(q_1, \dots, q_n, \phi(f), q) \mid (q_1, \dots, q_n, f, q) \in \delta\}.$$

By a trivial induction over the structure of the trees, it can be shown that

$$\phi(L(A)) = L(\phi(A)). \quad (5)$$

### 3 Position Functions

In this section, we define the position functions that are considered in the construction of the Bottom-Up automaton in the next sections. We show how to compute them and how they characterize the trees in the language denoted by a given expression.

Let  $E$  be a linear expression over a ranked alphabet  $\Sigma$  and  $f$  be a symbol  $\in \Sigma_k$ . The set  $\text{Root}(E)$ , subset of  $\Sigma$ , contains the roots of the trees in  $L(E)$ , *i.e.*

$$\text{Root}(E) = \{\text{root}(t) \mid t \in L(E)\}. \quad (6)$$

The set  $\text{Father}(E, f)$ , subset of  $\Sigma \times \mathbb{N}$ , contains a couple  $(g, i)$  if there exists a tree in  $L(E)$  with a node labeled by  $g$  the  $i$ -th child of is a node labeled by  $f$ :

$$\text{Father}(E, f) = \bigcup_{t \in L(E)} \text{father}(t, f). \quad (7)$$

*Example 1.* Let us consider the ranked alphabet defined by  $\Sigma_2 = \{f\}$ ,  $\Sigma_1 = \{g\}$ , and  $\Sigma_0 = \{a, b\}$ . Let  $E$  and  $\bar{E}$  be the expressions defined by

$$E = (f(a, a) + g(b))^* \cdot_b f(g(a), b), \quad \bar{E} = (f_1(a, a) + g_2(b))^* \cdot_b f_3(g_4(a), b).$$

Hence,

$$\begin{aligned} \text{Root}(\overline{E}) &= \{a, f_1, g_2\}, \\ \text{Father}(\overline{E}, f_1) &= \{(f_1, 1), (f_1, 2)\}, & \text{Father}(\overline{E}, a) &= \{(f_1, 1), (f_1, 2), (g_4, 1)\}, \\ \text{Father}(\overline{E}, g_2) &= \{(f_1, 1), (f_1, 2)\}, & \text{Father}(\overline{E}, b) &= \{(f_3, 2)\}, \\ \text{Father}(\overline{E}, f_3) &= \{(g_2, 1)\}, & \text{Father}(\overline{E}, g_4) &= \{f_3, 1\}. \end{aligned}$$

Let us show how to inductively compute these functions.

**Lemma 1.** *Let  $E$  be a linear expression over a ranked alphabet  $\Sigma$ . The set  $\text{Root}(E)$  is inductively computed as follows:*

$$\begin{aligned} \text{Root}(f(E_1, \dots, E_n)) &= \{f\}, \\ \text{Root}(E_1 + E_2) &= \text{Root}(E_1) \cup \text{Root}(E_2), \\ \text{Root}(E_1 \cdot_c E_2) &= \begin{cases} \text{Root}(E_1) \setminus \{c\} \cup \text{Root}(E_2) & \text{if } c \in L(E_1), \\ \text{Root}(E_1) & \text{otherwise,} \end{cases} \\ \text{Root}(E_1^{*c}) &= \text{Root}(E_1) \cup \{c\}, \end{aligned}$$

where  $E_1, \dots, E_n$  are  $n$  regular expressions over  $\Sigma$ ,  $f$  is a symbol in  $\Sigma_n$  and  $c$  is a symbol in  $\Sigma_0$ .

**Lemma 2.** *Let  $E$  be a linear expression and  $f$  be a symbol in  $\Sigma_k$ . The set  $\text{Father}(E, f)$  is inductively computed as follows:*

$$\begin{aligned} \text{Father}(g(E_1, \dots, E_n), f) &= \bigcup_{i \leq n} \text{Father}(E_i, f) \cup \{(g, i) \mid f \in \text{Root}(E_i)\}, \\ \text{Father}(E_1 + E_2, f) &= \text{Father}(E_1, f) \cup \text{Father}(E_2, f), \\ \text{Father}(E_1 \cdot_c E_2, f) &= (\text{Father}(E_1, f) \mid f \neq c) \cup \text{Father}(E_2, f) \\ &\quad \cup (\text{Father}(E_1, c) \mid f \in \text{Root}(E_2)) \\ \text{Father}(E_1^{*c}, f) &= \text{Father}(E_1, f) \cup (\text{Father}(E_1, c) \mid f \in \text{Root}(E_1)), \end{aligned}$$

where  $E_1, \dots, E_n$  are  $n$  regular expressions over  $\Sigma$ ,  $g$  is a symbol in  $\Sigma_n$  and  $c$  is a symbol in  $\Sigma_0$ .

*Proof (partial).* Let us consider the following cases.

- (1) Let us consider a tree  $t = t_1 \cdot_c L(E_2)$  with  $t_1 \in L(E_1)$ . By definition,  $t$  equals  $t_1$  where the occurrences of  $c$  have been replaced by some trees  $t_2$  in  $L(E_2)$ . Two cases may occur. **(a)** If  $c \neq f$ , then a predecessor of the symbol  $f$  in  $t$  can be a predecessor of the symbol  $f$  in a tree  $t_2$  in  $L(E_2)$ , a predecessor of the symbol  $f$  in  $t_1$ , or a predecessor of  $c$  in  $t_1$  if an occurrence of  $c$  in  $t_1$  has been replaced by a tree  $t_2$  in  $L(E_2)$  the root of which is  $f$ . **(b)** If  $c = f$ , since the occurrences of  $c$  have been replaced by some trees  $t_2$  of  $L(E_2)$ , a predecessor of the symbol  $c$  in  $t$  can be a predecessor of the symbol  $c$  in a tree  $t_2$  in  $L(E_2)$ , or a predecessor of  $c$  in  $t_1$  if an occurrence of  $c$  has been replaced by itself (and therefore if it appears in  $L(E_2)$ ). In both of these two cases, we conclude using Eqs. (3) and (6).

- (2) By definition,  $L(E_1^{*c}) = \bigcup_{k \geq 0} L(E_1)^{c,k}$ . Therefore, a tree  $t$  in  $L(E_1^{*c})$  is either  $c$  or a tree  $t_1$  in  $L(E_1)$  where the occurrences of  $c$  have been replaced by some trees  $t_2$  in  $L(E_1)^{c,k}$  for some integer  $k$ . Let us then proceed by recursion over  $k$ . If  $k = 1$ , a predecessor of  $f$  in  $t$  is a predecessor of  $f$  in  $t_1$ , a predecessor of  $f$  in a tree  $t_2$  in  $L(E_1)^{c,1}$  or a predecessor of  $c$  in  $t_1$  if an occurrence of  $c$  in  $t_1$  was substituted by a tree  $t_2$  in  $L(E_1)^{c,1}$  the root of which is  $f$ , *i.e.*

$$\text{Father}(E_1^{c,2}, f) = \text{Father}(E_1, f) \cup (\text{Father}(E_1, c) \mid f \in \text{Root}(E_1)).$$

By recursion over  $k$  and with the same reasoning, each recursion step adds  $\text{Father}(E_1, f)$  to the result of the previous step, and therefore

$$\text{Father}(E_1^{c,k}, f) = \text{Father}(E_1, f) \cup (\text{Father}(E_1, c) \mid f \in \text{Root}(E_1)).$$

□

Let us now show how these functions characterize, for a tree  $t$ , the membership of  $t$  in the language denoted by an expression.

**Definition 1.** Let  $E$  be a linear expression over a ranked alphabet  $\Sigma$  and  $t$  be a tree in  $T(\Sigma)$ . The property  $P(t)$  is the property defined by

$$\forall s = f(t_1, \dots, t_n) \prec t, \forall i \leq n, (f, i) \in \text{Father}(E, \text{root}(t_i)).$$

**Proposition 1.** Let  $E$  be a linear expression over a ranked alphabet  $\Sigma$  and  $t$  be a tree in  $T(\Sigma)$ . Then (1)  $t$  is in  $L(E)$  if and only if (2)  $\text{root}(t)$  is in  $\text{Root}(E)$  and  $P(t)$  is satisfied.

*Proof (partial).* Let us first notice that the proposition  $1 \Rightarrow 2$  is direct by definition of  $\text{Root}$  and  $\text{Father}$ . Let us show the second implication by induction over the structure of  $E$ . Hence, let us suppose that  $\text{root}(t)$  is in  $\text{Root}(E)$  and  $P(t)$  is satisfied.

- (1) Let us consider the case when  $E = E_1 \cdot_c E_2$ . Let us first suppose that  $\text{root}(t)$  is in  $\text{Root}(E_2)$ . Then  $c$  is in  $L(E_1)$  and  $P(t)$  is equivalent to

$$\forall s = f(t_1, \dots, t_n) \prec t, \forall i \leq n, (f, i) \in \text{Father}(E_2, \text{root}(t_i)).$$

By induction hypothesis  $t$  is in  $L(E_2)$  and therefore in  $L(E)$ .

Let us suppose now that  $\text{root}(t)$  is in  $\text{Root}(E_1)$ . Since  $E$  is linear, let us consider the subtrees  $t_2$  of  $t$  with only symbols of  $E_2$  and a symbol of  $E_1$  as a predecessor in  $t$ . Since  $P(t)$  holds, according to induction hypothesis and Lemma 2, each of these trees belongs to  $L(E_2)$ . Hence  $t$  belongs to  $t_1 \cdot_c L(E_2)$  where  $t_1$  is equal to  $t$  where the previously defined  $t_2$  trees are replaced by  $c$ . Once again, since  $P(t)$  holds and since  $\text{root}(t)$  is in  $\text{Root}(E_1)$ ,  $t_1$  belongs to  $L(E_1)$ .

In these two cases,  $t$  belongs to  $L(E)$ .

(2) Let us consider the case when  $E = E_1^{*c}$ . Let us proceed by induction over the structure of  $t$ . If  $t = c$ , the proposition holds from Lemmas 1 and 2. Following Lemma 2, each predecessor of a symbol  $f$  in  $t$  is a predecessor of  $f$  in  $E_1$  (case 1) or a predecessor of  $c$  in  $E_1$  (case 2). If all the predecessors of the symbols satisfy the case 1, then by induction hypothesis  $t$  belongs to  $L(E_1)$  and therefore to  $L(E)$ . Otherwise, we can consider (similarly to the catenation product case) the smallest subtrees  $t_2$  of  $t$  the root of which admits a predecessor in  $t$  which is a predecessor of  $c$  in  $E_1$ . By induction hypothesis, these trees belong to  $L(E_1)$ . And consequently  $t$  belongs to  $t' \cdots L(E_1)$  where  $t'$  is equal to  $t$  where the subtrees  $t_2$  have been substituted by  $c$ . Once again, by induction hypothesis,  $t'$  belongs to  $L(E_1^{*c})$ . As a direct consequence,  $t$  belongs to  $L(E)$ .  $\square$

## 4 Bottom-Up Position Automaton

In this section, we show how to compute a Bottom-Up automaton with a linear number of states from the position functions previously defined.

**Definition 2.** *The Bottom-Up position automaton  $\mathcal{P}_E$  of a linear expression  $E$  over a ranked alphabet  $\Sigma$  is the automaton  $(\Sigma, \text{Pos}(E), \text{Root}(E), \delta)$  defined by:*

$$((f_1, \dots, f_n), g, g) \in \delta \Leftrightarrow \forall i \leq n, (g, i) \in \text{Father}(E, f_i).$$

*Example 2.* The Bottom-Up position automaton  $(\text{Pos}(\overline{E}), \text{Pos}(\overline{E}), \text{Root}(\overline{E}), \delta)$  of the expression  $\overline{E}$  defined in Example 1 is defined as follows:

$$\begin{aligned} \text{Pos}(E) &= \{a, b, f_1, g_2, f_3, g_4\}, \text{Root}(\overline{E}) = \{a, f_1, g_2\}, \\ \delta &= \{(a, a), (b, b), ((a, a), f_1, f_1), ((a, f_1), f_1, f_1), ((a, g_2), f_1, f_1), ((f_1, a), f_1, f_1), \\ &\quad ((f_1, f_1), f_1, f_1), ((f_1, g_2), f_1, f_1), ((g_2, a), f_1, f_1), ((g_2, f_1), f_1, f_1), \\ &\quad ((g_2, g_2), f_1, f_1), (f_3, g_2, g_2), ((b, g_4), f_3, f_3), (a, g_4, g_4)\} \end{aligned}$$

Let us now show that the position automaton of  $E$  recognizes  $L(E)$ .

**Lemma 3.** *Let  $\mathcal{P}_E = (\Sigma, Q, Q_F, \delta)$  be the Bottom-Up position automaton of a linear expression  $E$  over a ranked alphabet  $\Sigma$ ,  $t$  be a tree in  $T_\Sigma$  and  $f$  be a symbol in  $\text{Pos}(E)$ . Then (1)  $f \in \Delta(t)$  if and only if (2)  $\text{root}(t) = f \wedge P(t)$ .*

*Proof.* Let us proceed by induction over the structure of  $t = f(t_1, \dots, t_n)$ . By definition,  $\Delta(t) = \delta(\Delta(t_1), \dots, \Delta(t_n), f)$ . For any state  $f_i$  in  $\Delta_i$ , it holds from the induction hypothesis that

$$f_i \in \Delta(t_i) \Leftrightarrow \text{root}(t_i) = f_i \wedge P(t_i). \quad (*)$$

Then, suppose that (1) holds (*i.e.*  $f \in \Delta(t)$ ). Equivalently, there exists by definition of  $\mathcal{P}_E$  a transition  $((f_1, \dots, f_n), f, f)$  in  $\delta$  such that  $f_i$  is in  $\Delta(t_i)$  for any integer  $i \leq n$ . Consequently,  $f$  is the root of  $t$ . Moreover, from the equivalence stated in Eq. (\*),  $\text{root}(t_i) = f_i$  and  $P(t_i)$  holds for any integer  $i \leq n$ . Finally and equivalently,  $P(t)$  holds as a consequence of Eq. (3). The reciprocal condition can be proved similarly since only equivalences are considered.  $\square$

As a direct consequence of Lemma 3 and Proposition 1.

**Proposition 2.** *The Bottom-Up position automaton of a linear expression  $E$  recognizes  $L(E)$ .*

The Bottom-Up position automaton of a (not necessarily linear) expression  $E$  can be obtained by first computing the Bottom-Up position automaton of its linearized expression  $\bar{E}$  and then by applying the alphabetical morphism  $h$ . As a direct consequence of Eq. (5).

**Proposition 3.** *The Bottom-Up position automaton of an expression  $E$  recognizes  $L(E)$ .*

## 5 Compressed Bottom-Up Position Automaton

In this section, we show that the structure of an expression allows us to factorize the transitions of a tree automaton by only considering the values of the Father function. The basic idea of the factorizations is to consider the cartesian product of sets. Imagine that a tree automaton contains four binary transitions  $(q_1, q_1, f, q_3)$ ,  $(q_1, q_2, f, q_3)$ ,  $(q_2, q_1, f, q_3)$  and  $(q_2, q_2, f, q_3)$ . These four transitions can be factorized as a *compressed transition*  $(\{q_1, q_2\}, \{q_1, q_2\}, f, q_3)$  using set of states instead of sets. The behavior of the original automaton can be simulated by considering the cartesian product of the origin states of the transition.

We first show how to encode such a notion of compressed automaton and how it can be used in order to solve the membership test.

**Definition 3.** *A compressed tree automaton over a ranked alphabet  $\Sigma$  is a 4-tuple  $(\Sigma, Q, Q_F, \delta)$  where  $Q$  is a set of states,  $Q_F \subset Q$  is the set of final states,  $\delta \subset (2^Q)^n \times \Sigma_n \times 2^Q$  is the set of compressed transitions that can be seen as a function from  $(2^Q)^k \times \Sigma_k$  to  $2^Q$  defined by*

$$(Q_1, \dots, Q_k, f, q) \in \delta \Leftrightarrow q \in \delta(Q_1, \dots, Q_k, f).$$

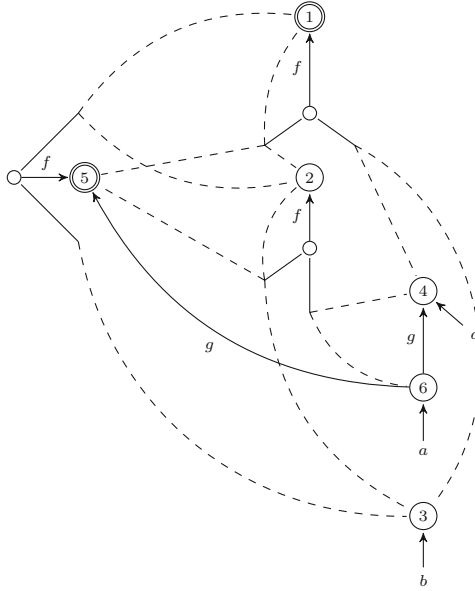
*Example 3.* Let us consider the compressed automaton  $A = (\Sigma, Q, Q_F, \delta)$  shown in Fig. 1. Its transitions are

$$\begin{aligned} \delta = \{ & (\{1, 2, 5\}, \{3, 4\}, f, 1), (\{2, 3, 5\}, \{4, 6\}, f, 2), \\ & (\{1, 2\}, \{3\}, f, 5), (\{6\}, g, 4), (\{6\}, g, 5), (a, 6), (a, 4), (b, 3) \}. \end{aligned}$$

The transition function  $\delta$  can be restricted to a function from  $Q^n \times \Sigma_n$  to  $2^Q$  (e.g. in order to simulate the behavior of an uncompressed automaton) by considering for a tuple  $(q_1, \dots, q_k)$  of states and a symbol  $f$  in  $\Sigma_k$  all the “active” transitions  $(Q_1, \dots, Q_k, f, q)$ , that are the transitions where  $q_i$  is in  $Q_i$  for  $i \leq k$ . More formally, for any states  $(q_1, \dots, q_k)$  in  $Q^k$ , for any symbol  $f$  in  $\Sigma_k$ ,

$$\delta(q_1, \dots, q_k, f) = \bigcup_{\substack{(Q_1, \dots, Q_k, f, q) \in \delta, \\ \forall i \leq k, q_i \in Q_i}} \{q\}. \quad (8)$$





**Fig. 1.** The compressed automaton  $A$ .

The transition set  $\delta$  can be extended to a function  $\Delta$  from  $T(\Sigma)$  to  $2^Q$  by inductively considering, for a tree  $f(t_1, \dots, t_k)$  the “active” transitions  $(Q_1, \dots, Q_k, f, q)$  once a subtree is read, that is when  $\Delta(q_i)$  and  $Q_i$  admits a common state for  $i \leq k$ . More formally, for any tree  $t = f(t_1, \dots, t_k)$  in  $T(\Sigma)$ ,

$$\Delta(t) = \bigcup_{\substack{(Q_1, \dots, Q_k, f, q) \in \delta, \\ \forall i \leq k, \Delta(t_i) \cap Q_i \neq \emptyset}} \{q\}.$$

As a direct consequence of the two previous equations,

$$\Delta(f(t_1, \dots, t_n)) = \bigcup_{(q_1, \dots, q_n) \in \Delta(t_1) \times \dots \times \Delta(t_n)} \delta(q_1, \dots, q_n, f). \tag{9}$$

The *language recognized by* a compressed automaton  $A = (\Sigma, Q, Q_F, \delta)$  is the subset  $L(A)$  of  $T(\Sigma)$  defined by

$$L(A) = \{t \in T(\Sigma) \mid \Delta(t) \cap Q_F \neq \emptyset\}.$$

*Example 4.* Let us consider the automaton of Fig. 1 and let us show that the tree  $t = f(f(b, a), g(a))$  belongs to  $L(A)$ . In order to do so, let us compute  $\Delta(t')$  for each subtree  $t'$  of  $t$ . First, by definition,

$$\Delta(a) = \{4, 6\}, \qquad \Delta(b) = \{3\}.$$

Since the only transition in  $\delta$  labeled by  $f$  containing 3 in its first origin set and 4 or 6 in its second is the transition  $(\{2, 3, 5\}, \{4, 6\}, f, 2)$ ,

$$\Delta(f(b, a)) = \{2\}.$$

Since the two transitions labeled by  $g$  are  $(\{6\}, g, 4)$  and  $(\{6\}, g, 5)$ ,

$$\Delta(g(a)) = \{4, 5\}.$$

Finally, there are two transitions labeled by  $f$  containing 2 in their first origin and 4 or 5 in its second:  $(\{2, 3, 5\}, \{4, 6\}, f, 2)$  and  $(\{1, 2, 5\}, \{3, 4\}, f, 1)$ . Therefore

$$\Delta(f(f(b, a), g(a))) = \{1, 2\}.$$

Finally, since 1 is a final state,  $t \in L(A)$ .

Let  $\phi$  be an alphabetical morphism between two alphabets  $\Sigma$  and  $\Sigma'$ . The *image* by  $\phi$  of a compressed automaton  $A = (\Sigma, Q, Q_F, \delta)$  is the compressed automaton  $\phi(A) = (\Sigma', Q, Q_F, \delta')$  where

$$\delta' = \{(Q_1, \dots, Q_n, \phi(f), q) \mid (Q_1, \dots, Q_n, f, q) \in \delta\}.$$

By a trivial induction over the structure of the trees, it can be shown that

$$L(\phi(A)) = \phi(L(A)). \tag{10}$$

Due to their inductive structure, regular expressions are naturally factorizing the structure of transitions of a Glushkov automaton. Let us now define the compressed position automaton of an expression.

**Definition 4.** *The compressed Bottom-Up position automaton  $\mathcal{C}(E)$  of a linear expression  $E$  is the automaton  $(\Sigma, \text{Pos}(E), \text{Root}(E), \delta)$  defined by*

$$\delta = \{(Q_1, \dots, Q_k, f, \{f\}) \mid Q_i = \{g \mid (f, i) \in \text{Father}(E, g)\}\}.$$

*Example 5.* Let us consider the expression  $\overline{E}$  defined in Example 1. The compressed automaton of  $\overline{E}$  is represented at Fig. 2.

As a direct consequence of Definition 4 and of Eq. (8),

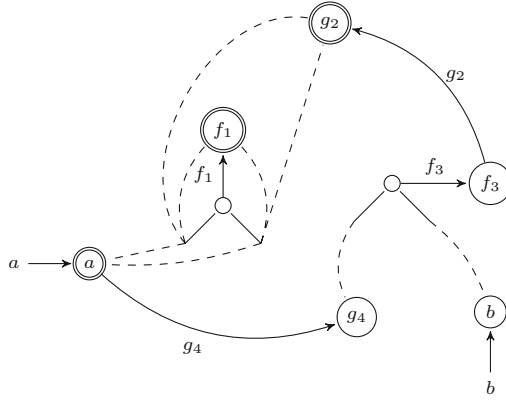
**Lemma 4.** *Let  $E$  be a linear expression over a ranked alphabet  $\Sigma$ . Let  $\mathcal{C}(E) = (\Sigma, Q, Q_F, \delta)$ . Then, for any states  $(q_1, \dots, q_n)$  in  $Q^n$ , for any symbol  $f$  in  $\Sigma_k$ ,*

$$\delta(q_1, \dots, q_n, f) = \{f\} \Leftrightarrow \forall i \leq n, (f, i) \in \text{Father}(E, q_i).$$

Consequently, considering Definition 2, Lemma 4 and Eq. (9),

**Proposition 4.** *Let  $E$  be a linear expression over a ranked alphabet  $\Sigma$ . Let  $\mathcal{P}_E = (-, -, -, \delta)$  and  $\mathcal{C}(E) = (-, -, -, \delta')$ . For any tree  $t$  in  $T(\Sigma)$ ,*

$$\Delta(t) = \Delta'(t).$$



**Fig. 2.** The compressed automata of the expression  $(f_1(a, a) + g_2(b))^* a \cdot_b f_3(g_4(a), b)$ .

Since the Bottom-Up position automaton of a linear expression  $E$  and its compressed version have the same states and the same final states,

**Corollary 1.** *The Glushkov automaton of an expression and its compact version recognize the same language.*

The compressed Bottom-Up position automaton of a (not necessarily linear) expression  $E$  can be obtained by first computing the compressed Bottom-Up position automaton of its linearized expression  $\bar{E}$  and then by applying the alphabetical morphism  $h$ . Therefore, considering Eq. (10),

**Proposition 5.** *The compressed Bottom-Up position automaton of a regular expression  $E$  recognizes  $L(E)$ .*

## 6 Web Application

The computation of the position functions and the Glushkov constructions have been implemented in a web application (made in Haskell, compiled into JavaScript using the REFLEX PLATFORM, represented with VIZ.JS) in order to help the reader to manipulate the notions. From a regular expression, it computes the classical Top-Down Glushkov defined in [8], and both the normal and the compressed versions of the Glushkov Bottom-Up automaton.

This web application can be found here [10]. As an example, the expression  $(f(a, a) + g(b))^* a \cdot_b f(g(a), b)$  of Example 1 can be defined from the literal input  $(f(a, a)+g(b))^* a \cdot_b f(g(a), b)$ .

## 7 Conclusion and Perspectives

In this paper, we have shown how to compute the Bottom-Up position automaton associated with a regular expression. This construction is relatively similar to

the classical one defined over a word expression [6]. We have also proposed a reduced version, the compressed Bottom-Up position automaton, that can be easily defined for word expressions too.

Since this construction is related to the classical one, one can wonder if all the studies involving Glushkov word automata can be extended to tree ones ([2–4, 12]). The classical Glushkov construction was also studied *via* its morphic links with other well-known constructions. The next step of our study is to extend Antimirov partial derivatives [1] in a Bottom-Up way too (in a different way from [7]), using the Bottom-Up quotient defined in [5].

## References

1. Antimirov, V.M.: Partial derivatives of regular expressions and finite automaton constructions. *Theor. Comput. Sci.* **155**(2), 291–319 (1996)
2. Bouchou, B., Duarte, D., Alves, M.H.F., Laurent, D., Musicante, M.A.: Schema evolution for XML: a consistency-preserving approach. In: Fiala, J., Koubek, V., Kratochvíl, J. (eds.) *MFCS 2004*. LNCS, vol. 3153, pp. 876–888. Springer, Heidelberg (2004). [https://doi.org/10.1007/978-3-540-28629-5\\_69](https://doi.org/10.1007/978-3-540-28629-5_69)
3. Brüggemann-Klein, A., Wood, D.: One-unambiguous regular languages. *Inf. Comput.* **140**(2), 229–253 (1998)
4. Caron, P., Ziadi, D.: Characterization of Glushkov automata. *Theor. Comput. Sci.* **233**(1–2), 75–90 (2000)
5. Champarnaud, J., Mignot, L., Sebti, N.O., Ziadi, D.: Bottom-up quotients for tree languages. *J. Autom. Lang. Comb.* **22**(4), 243–269 (2017)
6. Glushkov, V.M.: The abstract theory of automata. *Russ. Math. Surv.* **16**, 1–53 (1961)
7. Kuske, D., Meinecke, I.: Construction of tree automata from regular expressions. *RAIRO Theor. Inf. Appl.* **45**(3), 347–370 (2011)
8. Laugerotte, É., Sebti, N.O., Ziadi, D.: From regular tree expression to position tree automaton. In: Dediu, A.-H., Martín-Vide, C., Truthe, B. (eds.) *LATA 2013*. LNCS, vol. 7810, pp. 395–406. Springer, Heidelberg (2013). [https://doi.org/10.1007/978-3-642-37064-9\\_35](https://doi.org/10.1007/978-3-642-37064-9_35)
9. McNaughton, R.F., Yamada, H.: Regular expressions and state graphs for automata. *IEEE Trans. Electron. Comput.* **9**, 39–57 (1960)
10. Mignot, L.: Application: Glushkov tree automata. <http://ludovicmignot.free.fr/programmes/glushkovBotUp/index.html>. Accessed 27 Feb 2018
11. Mignot, L., Sebti, N.O., Ziadi, D.: Tree automata constructions from regular expressions: a comparative study. *Fundam. Inform.* **156**(1), 69–94 (2017)
12. Nicaud, C.: On the average size of Glushkov’s automata. In: Dediu, A.H., Ionescu, A.M., Martín-Vide, C. (eds.) *LATA 2009*. LNCS, vol. 5457, pp. 626–637. Springer, Heidelberg (2009). [https://doi.org/10.1007/978-3-642-00982-2\\_53](https://doi.org/10.1007/978-3-642-00982-2_53)