

CORRIGENDUM TO OUR PAPER: HOW EXPRESSIONS CAN CODE FOR AUTOMATA

SYLVAIN LOMBARDY¹ AND JACQUES SAKAROVITCH²

Abstract. We correct a mistake made in a previous paper in the construction of an automaton from a rational expression. We used there the definition of derivation of expression given by Antimirov, and this definition has to be further adapted for our purpose.

1991 Mathematics Subject Classification. 68Q45, 68Q70.

A DISTURBING EXAMPLE

In [7], we were considering the following problem: *Is it possible to build an algorithm Ω such that for any rational expression E computed from an automaton \mathcal{A} — i.e. $E = \Phi(\mathcal{A})$ where Φ is the state elimination method for instance — the following holds: $\mathcal{A} = \Omega(E)$?* We did not solve the problem completely, but we have identified two constructions that are good candidates to be the core components of such an algorithm Ω . The first one is the construction of an automaton $\Delta(E)$ from an expression E , the second one is the computation of the *minimal co-quotient* $\Upsilon(\mathcal{B})$ of an automaton \mathcal{B} . There is no problem with the minimal co-quotient but the definition we gave for $\Delta(E)$ was faulty. This can be observed for instance on the following example.

Let \mathcal{A}_1 be the automaton of Figure 1 (a), and let $E_1 = (a + b + 1)[a(a + b)]^*$ be the expression computed from \mathcal{A}_1 , which we write $E_1 = (a + b + 1)F_1$ with $F_1 = [a(a + b)]^*$; Figure 1 (b) shows $\Delta(E_1)$, whose co-quotient is not isomorphic to \mathcal{A}_1 .

As we shall see, it is not difficult to recover correct definitions and a true statement (Theorem 1.6) for the key result (Theorem 3.5) in the original paper.

Keywords and phrases: finite automata, regular expression, derivation of expressions, quotient of automata.

¹ IGM-LabInfo (UMR 8049), Université Paris-Est Marne-la-Vallée, 77454 Marne-la-Vallée Cedex 2, France; e-mail: lombardy@univ-mlv.fr.

² LTCI (UMR 5141), CNRS / ENST, 46 rue Barrault, 75634 Paris Cedex 13, France; e-mail: sakarovitch@enst.fr.

What proves to be more difficult, or at least much longer, is the complete and detailed proof, hopefully correct now, of Theorem 1.6.

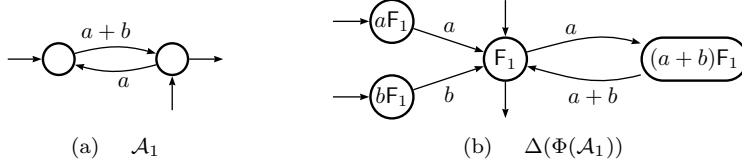


FIGURE 1. A counter-example to the statement in [7]

1. THE CORRECTED DEFINITIONS AND STATEMENT

In order to make this corrigendum as self-contained as possible, let us first recall the notion of derivation as defined by Antimirov. If E is a rational expression on A , we write $|E|$ for the language of A^* denoted by E and $c(E)$ for the Boolean whose value is 1 if $|E|$ contains the empty word and 0 otherwise; $c(E)$ is effectively computed by induction on the depth of E .

Definition 1.1 (Antimirov [2]). Let E be a rational expression on A and let a be a letter in A . The \mathbb{B} -derivative¹ of E with respect to a , denoted $\frac{\partial}{\partial a} E$, is a set of rational expressions on A , recursively defined by

$$\forall a, b \in A \quad \begin{aligned} \frac{\partial}{\partial a} 0 &= \frac{\partial}{\partial a} 1 = \emptyset, \\ \frac{\partial}{\partial a} b &= \begin{cases} \{1\} & \text{if } b = a, \\ \emptyset & \text{otherwise,} \end{cases} \end{aligned}$$

$$\frac{\partial}{\partial a} (E+F) = \frac{\partial}{\partial a} E \cup \frac{\partial}{\partial a} F, \quad (1)$$

$$\frac{\partial}{\partial a} (E \cdot F) = \left[\frac{\partial}{\partial a} E \right] \cdot F \cup c(E) \frac{\partial}{\partial a} F, \quad (2)$$

$$\frac{\partial}{\partial a} (E^*) = \left[\frac{\partial}{\partial a} E \right] \cdot E^*. \quad (3)$$

The induction implied by (1–3) should be interpreted while distributing derivation and product over union:

$$\frac{\partial}{\partial a} \left[\bigcup_{i \in I} E_i \right] = \bigcup_{i \in I} \frac{\partial}{\partial a} E_i, \quad \left[\bigcup_{i \in I} E_i \right] \cdot F = \bigcup_{i \in I} (E_i \cdot F).$$

¹We call it “ \mathbb{B} -derivative” and not simply “derivative” for two reasons. First in order to avoid confusion with the derivation defined by Brzozowski [3], and second because the formulae depend on the semiring of multiplicities and can be defined for other semirings (cf. [6]).

Definition 1.2. Let E be a rational expression on A and g a non empty word of A^* , i.e. $g = fa$ with a in A . The \mathbb{B} -derivative of E with respect to g , denoted $\frac{\partial}{\partial g} E$, is the set of rational expressions over A , recursively defined by formulae (1) – (3) and by

$$\forall f \in A^+, \forall a \in A \quad \frac{\partial}{\partial fa} E = \frac{\partial}{\partial a} \left(\frac{\partial}{\partial f} E \right) . \quad (4)$$

We call *true derived term* of E any rational expression which belongs to a set $\frac{\partial}{\partial g} E$ for some g in A^+ and *derived term* the expression E itself, or any true derived term; we write $D(E)$ for the set of derived terms of E .

It will be useful to distinguish the set of derived terms that are obtained as the result of a derivation by a word w : $D_w(E) = \{K \mid \exists u \in A^* \quad K \in \frac{\partial}{\partial uw} E\}$.

Definition 1.3. The *derived term automaton* of a rational expression E is the finite automaton $A(E)$ whose set of states is $D(E)$ and whose transitions are defined by:

- (i) if K and K' are derived terms of E and if a is a letter of A , (K, a, K') is a transition if and only if K' belongs to $\frac{\partial}{\partial a} K$;
- (ii) the initial state is E ;
- (iii) a derived term K is final if and only if $c(K) = 1$.

The essence of the derivation (by a letter), as defined by Antimirov and repeated here, is to “break” the expression into pieces when the operator at the upper level of the expression is “+”. The modification of the derivation we consider now consists in supposing that this breaking happens *spontaneously*, as if it were a derivation with respect to the empty word, before the first derivation by a letter, and after every such derivation. To that end, we define a new operation on rational expressions which we denote by $d()$ and which, roughly speaking, consists of decomposing an expression into a set of expressions whose left factor is not a sum.

Here lie the correction we make to our original paper. The definition of the derivation we apply to the expressions is modified in two ways: first, the breaking $d()$ is more ‘complete’ as it ‘goes through’ sets of expressions which contain 1; second, we replace the derivation by the breaking derivation (in the original paper, the breaking $d()$ was applied only once before any derivation).

The definition of $d()$ requires two further notations. If X is a *set of expressions*, then $[X]_p$ is the same set possibly without the expression 1: $[X]_p = X \setminus \{1\}$, and δ_X is the Boolean that takes the value 1 if the expression 1 belongs to X and 0 otherwise. For instance, $\delta_{[X]_p} = 0$ for any set X . If E is an expression, $d(E)$ is a set of expressions inductively defined by:

$$\begin{aligned} d(0) &= \{0\}, & d(1) &= \{1\}, & \forall a \in A \quad d(a) &= \{a\}, \\ d(E + F) &= d(E) \cup d(F), & d(E \cdot F) &= [d(E)]_p \cdot F \cup \delta_{d(E)} d(F), & d(E^*) &= \{E^*\}, \end{aligned}$$

and $d(E)$ is called the set of *initial broken (derived) terms* of E .

Definition 1.4. The *breaking derivation* of an expression E with respect to a letter a is defined as:

$$\frac{\partial_{\mathbf{b}}}{\partial a} E = d\left(\left[\frac{\partial}{\partial a} E\right]\right).$$

This breaking derivation is then extended to words by composition as in (4) and we call *true broken derived term* of E any rational expression which belongs to a set $\frac{\partial_{\mathbf{b}}}{\partial g} E$ for some non empty word g in A^+ :

$$\text{TBD}(E) = \{L \mid \exists g \in A^+ \quad L \in \frac{\partial_{\mathbf{b}}}{\partial g} E\} .$$

We call *broken derived term* of E any element of the union of $d(E)$ and $\text{TBD}(E)$:

$$\text{BD}(E) = d(E) \cup \text{TBD}(E) .$$

It is easy to check that

$$\forall f \in A^* \quad \frac{\partial_{\mathbf{b}}}{\partial f} E = d\left(\left[\frac{\partial}{\partial f} E\right]\right) ,$$

which amounts to say that the broken derived terms are obtained by “breaking” the derived terms:

$$\text{BD}(E) = \bigcup_{K \in D(E)} d(K) . \quad (5)$$

The *broken derived term automaton* $\Delta'(E)$ is then defined as was the derived term automaton by Antimirov. Its set of states is $\text{BD}(E)$, the set of initial states is $d(E)$, a state K is final if and only if $c(K) = 1$ and for a in A , and K, K' in $\text{BD}(E)$, (K, a, K') is a transition if and only if K' belongs to $\frac{\partial_{\mathbf{b}}}{\partial a} K$.

The automaton $\Delta'(E)$ recognizes the language denoted by E (cf. [6, Th. 9]).² More precisely, it holds

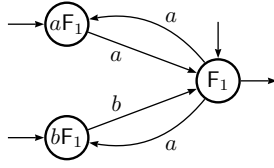
Property 1.5. For every K in $\text{BD}(E)$, the future of K in $\Delta'(E)$ is equal to $|K|$.

The corrected version for Theorem 3.5 of [7] is then obtained by replacing *derived term* by *broken derived term*.

Theorem 1.6. Let \mathcal{A} be a co-deterministic automaton and $E = \Phi(\mathcal{A})$ a rational expression computed from \mathcal{A} by the state elimination method. Then, the broken derived term automaton $\Delta'(E)$ of E is co-deterministic.

Figure 2 shows $\Delta'(E_1)$, a co-deterministic automaton as expected, and whose minimal co-quotient is \mathcal{A}_1 , as desired.

²We use that notation $\Delta'()$ here, in this corrigendum, in order to have a different notation from the one used in [7]; it might be the case that in further publications, we shall use $A()$ which is more natural.

FIGURE 2. $\Delta'(\mathbf{E}_1)$

In [7], we have given a number of examples where the mapping $\Upsilon \circ \Delta'$ applied to an expression that is computed from a co-minimal automaton \mathcal{A} yields \mathcal{A} itself. But the only case where we have been able to prove the property is when \mathcal{A} is co-deterministic. Moreover, this result gives the possibility of dealing with the general case by means of tagging. This is the reason why we consider Theorem 1.6 as the key statement of [7].

The introduction of the breaking operation is necessary to deal with automata with several initial states. Nevertheless, the core of the proof of Theorem 1.6 is Proposition 3.5 which establishes that the derived term automaton of an expression computed from a *normalized* co-deterministic automaton is co-deterministic, that is, a property of Antimirov derivation. A remarkable feature of this result is that its proof goes by induction on the number of states of the normalized automaton we start from. The main ingredients of this proof are first Lemma 2.1 that describes how the (Antimirov) derivation translates under a *continuous rational substitution* and then on Proposition 2.8 which states that every (broken) derived term is contained in the *future* of some state of the automaton on which the expression has been computed. The transfer from the hypothesis of Proposition 3.5 (normalized co-deterministic automata) to the one of Theorem 1.6 (co-deterministic automata) is made possible by Lemma 2.7.

2. PRELIMINARY TO THE PROOF

We first establish two properties of derivation that will be used in the sequel and whose scope is wider than the statement alone. For the ease of the proof, we also rather work with *normalized* automata. We then show, in two steps, that we can make this assumption without loss of generality.

2.1. A SUBSTITUTION LEMMA

We call *rational substitution* a map $\varphi: B \rightarrow \text{Rat}E A^*$ which is extended to $\varphi: \text{Rat}E B^* \rightarrow \text{Rat}E A^*$ by replacing every atom of a rational expression over B^* by its image under φ . Substitutions are consistent with rational operations in the sense that for every $E, F \in \text{Rat}E B^*$ it holds:

$$\varphi(E + F) = \varphi(E) + \varphi(F), \quad \varphi(E \cdot F) = \varphi(E) \cdot \varphi(F) \quad \text{and} \quad (\varphi(E))^* = \varphi((E)^*).$$

Naturally, φ also induces a substitution from B^* into A^* and it holds $|\varphi(\mathbf{E})| = \varphi(|\mathbf{E}|)$. A substitution φ is *continuous* if $c(\varphi(b))$ is null for every b in B . If φ is *continuous*, it holds:

$$\forall \mathbf{E} \in \text{RatE } B^* \quad c(\varphi(\mathbf{E})) = c(\mathbf{E}) . \quad (6)$$

Lemma 2.1. *Let $\varphi: \text{RatE } B^* \rightarrow \text{RatE } A^*$ be a continuous rational substitution. Then*

$$\forall a \in A \quad \frac{\partial}{\partial a} \varphi(\mathbf{E}) = \bigcup_{b \in B} \left[\frac{\partial}{\partial a} \varphi(b) \right] \varphi \left(\frac{\partial}{\partial b} \mathbf{E} \right) . \quad (7)$$

Proof. By induction on the depth of \mathbf{E} . If $\mathbf{E} = 0$ or 1 , both sides of (7) are zero; if $\mathbf{E} = b$, the right handside of (7) reduces to $\frac{\partial}{\partial a} \varphi(b)$ which is the left handside, and the base of the induction is established. The following three sequences of equalities give the three possible induction steps.

$$\begin{aligned} \frac{\partial}{\partial a} \varphi(\mathbf{F} + \mathbf{G}) &= \frac{\partial}{\partial a} [\varphi(\mathbf{F}) + \varphi(\mathbf{G})] = \frac{\partial}{\partial a} \varphi(\mathbf{F}) \cup \frac{\partial}{\partial a} \varphi(\mathbf{G}) \\ &= \bigcup_{b \in B} \left[\frac{\partial}{\partial a} \varphi(b) \right] \varphi \left(\frac{\partial}{\partial b} \mathbf{F} \right) \cup \bigcup_{b \in B} \left[\frac{\partial}{\partial a} \varphi(b) \right] \varphi \left(\frac{\partial}{\partial b} \mathbf{G} \right) \\ &= \bigcup_{b \in B} \left[\frac{\partial}{\partial a} \varphi(b) \right] \left[\varphi \left(\frac{\partial}{\partial b} \mathbf{F} \right) \cup \varphi \left(\frac{\partial}{\partial b} \mathbf{G} \right) \right] \\ &= \bigcup_{b \in B} \left[\frac{\partial}{\partial a} \varphi(b) \right] \varphi \left(\frac{\partial}{\partial b} (\mathbf{F} + \mathbf{G}) \right) \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial a} \varphi(\mathbf{F} \cdot \mathbf{G}) &= \frac{\partial}{\partial a} [\varphi(\mathbf{F}) \cdot \varphi(\mathbf{G})] = \left[\frac{\partial}{\partial a} \varphi(\mathbf{F}) \right] \varphi(\mathbf{G}) \cup c(\varphi(\mathbf{F})) \frac{\partial}{\partial a} \varphi(\mathbf{G}) \\ &= \left[\bigcup_{b \in B} \left[\frac{\partial}{\partial a} \varphi(b) \right] \varphi \left(\frac{\partial}{\partial b} \mathbf{F} \right) \right] \cdot \varphi(\mathbf{G}) \cup c(\mathbf{F}) \bigcup_{b \in B} \left[\frac{\partial}{\partial a} \varphi(b) \right] \varphi \left(\frac{\partial}{\partial b} \mathbf{G} \right) \\ &= \bigcup_{b \in B} \left[\frac{\partial}{\partial a} \varphi(b) \right] \varphi \left(\frac{\partial}{\partial b} (\mathbf{F} \cdot \mathbf{G}) \right) \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial a} \varphi(\mathbf{F}^*) &= \frac{\partial}{\partial a} (\varphi(\mathbf{F}))^* = \left[\frac{\partial}{\partial a} \varphi(\mathbf{F}) \right] \cdot (\varphi(\mathbf{F}))^* \\ &= \left[\bigcup_{b \in B} \left[\frac{\partial}{\partial a} \varphi(b) \right] \varphi \left(\frac{\partial}{\partial b} \mathbf{F} \right) \right] \cdot (\varphi(\mathbf{F}))^* \\ &= \bigcup_{b \in B} \left[\frac{\partial}{\partial a} \varphi(b) \right] \left[\varphi \left(\frac{\partial}{\partial b} \mathbf{F} \right) \cdot (\varphi(\mathbf{F}))^* \right] = \bigcup_{b \in B} \left[\frac{\partial}{\partial a} \varphi(b) \right] \varphi \left(\frac{\partial}{\partial b} \mathbf{F}^* \right) \end{aligned}$$

□

Note that Lemma 2.1 refers to the *derivation* and not to the breaking derivation; it turns out that we shall make use of the two notions in parallel.

2.2. NORMALIZED AUTOMATA

If \mathcal{A} is an automaton over A , we write $|\mathcal{A}|$ for the language of A^* denoted by \mathcal{A} .

Definition 2.2. An automaton is *standard* if it has exactly one initial state with no incoming transition.

To any automaton $\mathcal{A} = \langle Q, A, E, I, T \rangle$ we associate the standard automaton $\mathcal{A}_\$ = \langle Q \cup \{i\}, A \cup \{\$, \}, F, \{i\}, T \rangle$, with

$$F = E \cup \{(i, \$, p) \mid p \in I\}$$

and where $\$$ is not in A and i does not belong to Q . Note that \mathcal{A} is co-deterministic if, and only if, so is $\mathcal{A}_\$$. Clearly $|\mathcal{A}_\$| = \$|\mathcal{A}|$. This equality indeed generalizes to the broken derived terms in the following way.

Lemma 2.3. *Let \mathcal{A} be an automaton over A , $\mathcal{A}_\$$ the standard automaton defined as above, $E = \Phi_\omega(\mathcal{A})$ and $E_\$ = \Phi_\omega(\mathcal{A}_\$)$ the expressions obtained by the state elimination method (with respect to the same order ω on the states of \mathcal{A}). It then holds:*

$$d(E) = \frac{\partial_b}{\partial \$} E_\$ \quad \text{and} \quad \text{BD}(E) = \text{TBD}(E_\$) .$$

Definition 2.4. An automaton is *co-standard* if it has exactly one final with no outgoing transition, that is, if its transposition is standard.

To any automaton $\mathcal{A} = \langle Q, A, E, I, T \rangle$ we associate the co-standard automaton $\mathcal{A}_\mathcal{L} = \langle Q \cup \{t\}, A \cup \{\mathcal{L}\}, F, I, \{t\} \rangle$, with

$$F = E \cup \{(p, \mathcal{L}, t) \mid p \in T\}$$

and where \mathcal{L} is not in A and t does not belong to Q . Note that \mathcal{A} is co-deterministic if, and only if, so is $\mathcal{A}_\mathcal{L}$. Clearly $|\mathcal{A}_\mathcal{L}| = |\mathcal{A}| \mathcal{L}$. This equality indeed generalizes to the broken derived terms in the following way.

Lemma 2.5. *Let \mathcal{A} be an automaton over A , $\mathcal{A}_\mathcal{L}$ the co-standard automaton defined as above, $E = \Phi_\omega(\mathcal{A})$ and $E_\mathcal{L} = \Phi_\omega(\mathcal{A}_\mathcal{L})$. The projection $\pi: (A \cup \{\mathcal{L}\})^* \longrightarrow A^*$ induces a bijection between $\text{BD}(E_\mathcal{L}) \setminus \{1\}$ and $\text{BD}(E)$ and in particular between $d(E_\mathcal{L})$ and $d(E)$.*

Proof. Let $\mathcal{F}_\mathcal{L}$ be the family of expressions in $\text{Rat}E(A \cup \mathcal{L})$ that denote rational subsets of A^* followed by \mathcal{L} :

$$\mathcal{F}_\mathcal{L} = \{E \in \text{Rat}E(A \cup \mathcal{L}) \mid |E| = L \mathcal{L}, L \in \text{Rat} A^*\} .$$

Since any word in the language denoted by an element of $\mathcal{F}_\mathcal{L}$ does not contain the symbol \mathcal{L} but at the end, any element of $\mathcal{F}_\mathcal{L}$ is of one of the following forms: \mathcal{L} , $E_\mathcal{L} + F_\mathcal{L}$, or $G F_\mathcal{L}$, where $E_\mathcal{L}$ and $F_\mathcal{L}$ are in $\mathcal{F}_\mathcal{L}$ and G is in $\text{Rat}E A$.

As the algorithm Φ acts symbolically on the labels of the transitions, $\pi(\mathbf{E}_{\mathcal{L}}) = \mathbf{E}$. By induction on the depth of the expression, we first show that $\pi(\mathbf{d}(\mathbf{E}_{\mathcal{L}})) = \mathbf{d}(\pi(\mathbf{E}_{\mathcal{L}})) = \mathbf{d}(\mathbf{E})$. The interesting case is when $\mathbf{E}_{\mathcal{L}} = \mathbf{G}\mathbf{F}_{\mathcal{L}}$:

$$\pi(\mathbf{d}(\mathbf{G}\mathbf{F}_{\mathcal{L}})) = \pi(\mathbf{d}(\mathbf{G}) \mathbf{F}_{\mathcal{L}}) \cup \delta_{\mathbf{d}(\mathbf{G})} \pi(\mathbf{d}(\mathbf{F}_{\mathcal{L}})) = \mathbf{d}(\mathbf{G}) \mathbf{F} \cup \delta_{\mathbf{d}(\mathbf{G})} \mathbf{d}(\mathbf{F}) = \mathbf{d}(\mathbf{G}\mathbf{F}) \quad .$$

In the same way, for any $\mathbf{H}_{\mathcal{L}}$ in $\mathcal{F}_{\mathcal{L}}$, and for any letter a in A , it holds:

$$\pi\left(\frac{\partial_{\mathbf{b}}}{\partial a} \mathbf{H}_{\mathcal{L}}\right) = \frac{\partial_{\mathbf{b}}}{\partial a} \pi(\mathbf{H}_{\mathcal{L}}) \quad . \quad (8)$$

The statement is then shown by induction on the length of the derivation. Every (broken) derived term of $\mathbf{E}_{\mathcal{L}}$ different from 1 is in $\mathcal{F}_{\mathcal{L}}$. Let $\mathbf{K}_{\mathcal{L}}$ be a derived term in $\text{BD}(\mathbf{E}_{\mathcal{L}})$ and \mathbf{K} its projection, in $\text{BD}(\mathbf{E})$ by induction hypothesis. The derivation of $\mathbf{K}_{\mathcal{L}}$ with respect to \mathcal{L} gives either 1, or the empty set. The derivation of $\mathbf{K}_{\mathcal{L}}$ with respect to any letter a in A is, by (8), in bijection with the derivation of \mathbf{K} with respect to a . \square

Putting the standardisation and co-standardisation together, we get the following definition and statement.

Definition 2.6. An automaton is *normalized* if and only if it has exactly one initial state and one final state such that there is no incoming transition on the initial state and no outgoing transition from the final state.

To any automaton $\mathcal{A} = \langle Q, A, E, I, T \rangle$ we associate the normalized automaton $\mathcal{A}_n = \langle Q \cup \{i, t\}, A \cup \{\$, \mathcal{L}\}, F, \{i\}, \{t\} \rangle$, with

$$F = E \cup \{(i, \$, p) \mid p \in I\} \cup \{(p, \mathcal{L}, t) \mid p \in T\}$$

and where $\$$ and \mathcal{L} are not in A and i and t do not belong to Q . Note that \mathcal{A} is co-deterministic if, and only if, so is \mathcal{A}_n . Clearly $|\mathcal{A}_n| = \$|\mathcal{A}|\mathcal{L}$. This equality indeed generalizes to the broken derived terms in the following way.

Lemma 2.7. *Let \mathcal{A} be an automaton over A , \mathcal{A}_n the normalized automaton defined as above, $\mathbf{E} = \Phi_{\omega}(\mathcal{A})$ and $\mathbf{E}_n = \Phi_{\omega}(\mathcal{A}_n)$. The projection $\pi: (A \cup \{\mathcal{L}\})^* \rightarrow A^*$ induces a bijection between $\frac{\partial_{\mathbf{b}}}{\partial \$} \mathbf{E}_n$ and $\mathbf{d}(\mathbf{E})$ on one hand and between $\text{TBD}(\mathbf{E}_n) \setminus \{1\}$ and $\text{BD}(\mathbf{E})$ on the other hand. \square*

2.3. FUTURE OF STATES AND DERIVED TERMS

Let $\mathcal{A} = \langle Q, A, E, I, T \rangle$ be an automaton over A^* . For each state q of \mathcal{A} , the *past* of q (in \mathcal{A}) is the set of labels of computation which go from an initial state of \mathcal{A} to q , and we write it $\text{Past}_{\mathcal{A}}(q)$; the *future* of q (in \mathcal{A}) is the set of labels of computations that go from q to a final state of \mathcal{A} and we write it $\text{Fut}_{\mathcal{A}}(q)$:

$$\text{Past}_{\mathcal{A}}(q) = \{w \in A^* \mid \exists i \in I \quad i \xrightarrow{\mathcal{A}} w \rightarrow q\}, \quad \text{Fut}_{\mathcal{A}}(q) = \{w \in A^* \mid \exists t \in T \quad q \xrightarrow{\mathcal{A}} w \rightarrow t\}.$$

An automaton \mathcal{A} is *deterministic* if, and only if, the pasts of states are pairwise disjoint, and dually, \mathcal{A} is *co-deterministic* if, and only if, the futures of states are pairwise disjoint.

Proposition 2.8. *Let \mathcal{A} be a standard automaton with set of states Q and $\mathbf{E} = \Phi(\mathcal{A})$. For every derived term \mathbf{K} of \mathbf{E} there exists a state q in Q such that $|\mathbf{K}|$ is included in the future of q .*

Proof. We first transform $\mathcal{A} = \langle Q, A, E, I, T \rangle$ into an automaton $\widehat{\mathcal{A}} = \langle Q, A \times Q, \widehat{E}, I, T \rangle$ by relabelling every transition (p, a, q) in E as (p, a_q, q) . The projection π maps $A \times Q$ onto A and is extended to a map from $\text{RatE}(A \times Q)$ onto $\text{RatE } A$.

As π is a continuous (rational) substitution, with the property that for any a, b in A , (and any q in Q), $\frac{\partial}{\partial a} \pi(b_q)$ is equal to 1 or 0 according to whether a is equal to b or not, Lemma 2.1 implies that for any a in A and for any expression \mathbf{F} in $\text{RatE}(A \times Q)$ it holds

$$\frac{\partial}{\partial a} \pi(\mathbf{F}) = \bigcup_{q \in Q} \left[\pi \left(\frac{\partial}{\partial a_q} \mathbf{F} \right) \right] . \quad (9)$$

Let $\widehat{\mathbf{E}} = \Phi(\widehat{\mathcal{A}})$. As the algorithm Φ acts symbolically on the labels of the transitions to build the expression, it holds: $\pi(\widehat{\mathbf{E}}) = \mathbf{E}$. Then, by iteration of (9), for any \mathbf{K} in $\text{D}(\widehat{\mathbf{E}})$ there exists \mathbf{K}' in $\text{D}(\widehat{\mathbf{E}})$ such that $\mathbf{K} = \pi(\mathbf{K}')$. Thus, there exists w in $(A \times Q)^*$ such that \mathbf{K}' is in $\frac{\partial}{\partial w} \widehat{\mathbf{E}}$ and then

$$|\mathbf{K}'| \subseteq \left| \frac{\partial}{\partial w} \widehat{\mathbf{E}} \right| = w^{-1} |\widehat{\mathcal{A}}| .$$

If $\mathbf{K}' = \widehat{\mathbf{E}}$, then $\mathbf{K} = \mathbf{E}$ and $|\mathbf{K}'|$ is (in) the future of the (unique) initial state. Otherwise, w is a non empty word; let a_q is the last letter of w , $w^{-1} |\widehat{\mathcal{A}}|$ is contained in $\text{Fut}_{\widehat{\mathcal{A}}}(q)$. The projection by π gives

$$\text{Fut}_{\mathcal{A}}(q) = \pi(\text{Fut}_{\widehat{\mathcal{A}}}(q)) \quad \text{and thus} \quad |\mathbf{K}| \subseteq \text{Fut}_{\mathcal{A}}(q) .$$

□

The same property holds then for broken derived terms.

Corollary 2.9. *Let \mathcal{A} be an automaton with set of states Q and $\mathbf{E} = \Phi_\omega(\mathcal{A})$. For every broken derived term \mathbf{K} of \mathbf{E} there exists a state q in Q such that $|\mathbf{K}|$ is included in the future of q .*

Proof. Let \mathcal{A}_\S be the standard automaton associated with \mathcal{A} as in Section 2.2 and let $\mathbf{E}_\S = \Phi_\omega(\mathcal{A}_\S)$. For every \mathbf{K}_\S in $\text{BD}(\mathbf{E}_\S)$ there exists \mathbf{H}_\S in $\text{D}(\mathbf{E}_\S)$ such that \mathbf{K}_\S is in $\text{d}(\mathbf{H}_\S)$. Moreover, there exists, by Proposition 2.8, a state q in Q such that

$|H_{\S}| \subseteq \text{Fut}_{\mathcal{A}_{\S}}(q)$ hence $|K_{\S}| \subseteq \text{Fut}_{\mathcal{A}_{\S}}(q)$. By Lemma 2.3, $\text{TBD}(E_{\S}) = \text{BD}(E)$, therefore $|K| \subseteq \text{Fut}_{\mathcal{A}}(q)$. \square

As the futures of states are pairwise disjoint in a co-deterministic automaton, we then can state the following.

Corollary 2.10. *Let \mathcal{A} be a co-deterministic normalized automaton with set of states Q and $E = \Phi(\mathcal{A})$. For every K of $D(E)$, there exists a unique q in Q such that $|K| \subseteq \text{Fut}_{\mathcal{A}}(q)$; for every L of $\text{BD}(E)$, there exists a unique q in Q such that $|L| \subseteq \text{Fut}_{\mathcal{A}}(q)$. \square*

3. PROOF OF THEOREM 1.6

The proof goes by induction on the number of states of the automaton \mathcal{A} as sketched in Figure 3.

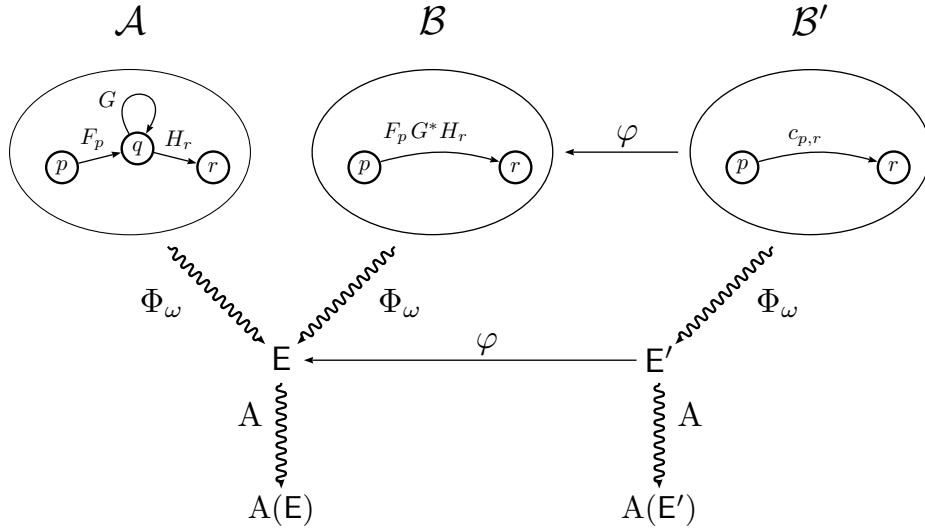


FIGURE 3. Diagram for an induction step in the proof of Theorem 1.6

Let $E = \Phi_\omega(\mathcal{A})$ the rational expression computed by the state elimination method following the order ω on the states of \mathcal{A} and let q be the smallest state with respect to ω . For every predecessor p of q , let F_p be the set of letters a such that (p, a, q) is a transition of \mathcal{A} and F_p the expression ‘sum of letters in F_p ’. Let R be the set of successors of q in \mathcal{A} ; for every r in R , let H_r be the set of letters a such that (q, a, r) is a transition of \mathcal{A} and H_r the expression ‘sum of letters in H_r ’. Let G be the set of letters a such that (q, a, q) is a transition of \mathcal{A} and G the expression ‘sum of letters in G ’ (G may be empty and G null).

Let \mathcal{B} be the *generalized* automaton obtained from \mathcal{A} by elimination of q — by the state elimination method. For every p and every r as above, a transition

from p to r labelled by $F_p G^* H_r$ (*resp.* $F_p H_r$ if G is null) is thus created.³ Let \mathcal{B}' be the (normalized) automaton obtained from \mathcal{B} by replacing every label of the form $F_p G^* H_r$ (or $F_p H_r$) by a fresh letter $c_{p,r}$. Let C be the alphabet of these fresh letters and $A' = A \cup C$. Let φ be the (continuous) substitution from $\text{RatE } A'^*$ to $\text{RatE } A^*$ which maps every letter of A onto itself and every letter $c_{p,r}$ onto $F_p G^* H_r$ (or $F_p H_r$ if G is null).

Since the construction of \mathcal{B} is the first step of $\Phi_\omega(\cdot)$, then $\mathbf{E} = \Phi_\omega(\mathcal{B})$ and let $\mathbf{E}' = \Phi_\omega(\mathcal{B}')$. As the algorithm Φ_ω acts symbolically on the labels of the transitions, $\varphi(\mathbf{E}') = \mathbf{E}$.

If \mathcal{A} is co-deterministic, so is \mathcal{B}' and $\Delta'(\mathbf{E}')$ is co-deterministic by induction hypothesis. The proof amounts to transfer the properties from $\Delta'(\mathbf{E}')$ to $\Delta'(\mathbf{E})$ via the substitution φ .

3.1. PREPARATION FOR THE INDUCTION

From now on — and until after Corollary 3.6 — \mathcal{A} , and thus \mathcal{B} and \mathcal{B}' , are *normalized* automaton. The state elimination method does not eliminate neither the initial nor the final state and thus every state is treated in the same way. With the above notation, it first holds:

Property 3.1. *For every state s , $s \neq q$, $\text{Fut}_{\mathcal{A}}(s) = \varphi(\text{Fut}_{\mathcal{B}'}(s))$. □*

The remainder of this subsection is devoted to the description of the image of the broken derived terms under φ . From Lemma 2.1, for every K' in $\text{RatE } A'^*$, and every letter a in A , it holds

$$\frac{\partial}{\partial a} \varphi(K') = \varphi \left(\frac{\partial}{\partial a} K' \right) \cup \bigcup_{c \in C} \left[\frac{\partial}{\partial a} \varphi(c) \right] \varphi \left(\frac{\partial}{\partial c} K' \right) .$$

For every c in C , there exists p and r such that $\varphi(c) = F_p G^* H_r$ (*resp.* $F_p H_r$). Hence, either a does not belong to any F_p and $\frac{\partial}{\partial a} \varphi(K') = \varphi \left(\frac{\partial}{\partial a} K' \right)$ or it belongs to F_p for a unique p and (recall that R is the set of successors of q in \mathcal{A}):

$$\frac{\partial}{\partial a} \varphi(K') = \begin{cases} \varphi \left(\frac{\partial}{\partial a} K' \right) \cup \bigcup_{r \in R} G^* H_r \varphi \left(\frac{\partial}{\partial c_{p,r}} K' \right) & \text{if } G \text{ is not empty,} \\ \varphi \left(\frac{\partial}{\partial a} K' \right) \cup \bigcup_{r \in R} H_r \varphi \left(\frac{\partial}{\partial c_{p,r}} K' \right) & \text{otherwise.} \end{cases} \quad (10)$$

As G and H_r are sums of letters, for every c in C the derivation of $G^* H_r \varphi \left(\frac{\partial}{\partial c} K' \right)$ (*resp.* $H_r \varphi \left(\frac{\partial}{\partial c} K' \right)$) with respect to any letter gives \emptyset , $G^* H_r \varphi \left(\frac{\partial}{\partial c} K' \right)$ or

³The automaton \mathcal{B} is labelled by rational expressions, which is the reason why it is called ‘generalized’.

$\varphi\left(\frac{\partial}{\partial c} K'\right)$. Hence, by a straightforward induction, it holds :

$$D(E) = \varphi(D(E')) \cup \bigcup_{c_{p,r} \in C} \left[\bigcup_{a \in F_p} \frac{\partial}{\partial a} \varphi(c_{p,r}) \right] \varphi(D_{c_{p,r}}(E')) . \quad (11)$$

By applying $d()$ to (10), it comes $\frac{\partial_b}{\partial a} \varphi(K') = \varphi\left(\frac{\partial_b}{\partial a} K'\right)$ in the former case, and in the latter:

$$\frac{\partial_b}{\partial a} \varphi(K') = \begin{cases} \varphi\left(\frac{\partial_b}{\partial a} K'\right) \cup \bigcup_{r \in R} G^* H_r \varphi\left(\frac{\partial}{\partial c_{p,r}} K'\right) & \text{if } G \text{ is not empty,} \\ \varphi\left(\frac{\partial_b}{\partial a} K'\right) \cup \bigcup_{r \in R} \bigcup_{b \in H_r} b \varphi\left(\frac{\partial}{\partial c_{p,r}} K'\right) & \text{otherwise.} \end{cases} \quad (12)$$

As above, and for every c in C the *breaking* derivation of $G^* H_r \varphi\left(\frac{\partial}{\partial c} K'\right)$ (*resp.* $b \varphi\left(\frac{\partial}{\partial c} K'\right)$) with respect to any letter gives \emptyset , $G^* H_r \varphi\left(\frac{\partial}{\partial c} K'\right)$ or $\varphi\left(\frac{\partial_b}{\partial c} K'\right)$. Hence, by a straightforward induction, it holds :

$$BD(E) = \varphi(BD(E')) \cup \bigcup_{c_{p,r} \in C} \left[\bigcup_{a \in F_p} \frac{\partial_b}{\partial a} \varphi(c_{p,r}) \right] \varphi(D_{c_{p,r}}(E')) . \quad (13)$$

3.2. PROOF OF THEOREM 1.6

The comparison between (11) and (13) is quite instructive. The presence in the right handside of both equations of a *non breaking derivation* makes it understandable that an induction proof will run more smoothly on the derived term automata. This is done with Proposition 3.5. On the other hand, when it comes to compare the derived terms of an expression obtained from an automaton \mathcal{A} and those obtained from its normalization \mathcal{A}_s , the *broken derived terms* are to be used (Lemma 2.7). Corollary 3.6 will fill the gap and Theorem 1.6 will follow.

We keep the notation of the preceding subsection and we suppose from now on that \mathcal{A} is (normalized) *co-deterministic*. This implies that \mathcal{B}' is co-deterministic and moreover that for every p, G and F_p are disjoint; likewise, if p and p' are two distinct predecessors of q , F_p and $F_{p'}$ are disjoint. For every successor r of q , let I_r be the set of labels of incoming transitions of r that does not come from q ; then I_r and H_r are disjoint. For every state s of \mathcal{B}' (therefore different from q), let J_s be the set of labels of incoming transitions of s . It then holds:

Property 3.2. *For every state s in \mathcal{B}' , $\varphi(J_s)$ is a suffix code in A^* and the images of the elements of J_s by φ are pairwise disjoint.*

Lemma 3.3. *For every state s in \mathcal{A} , $s \neq q$, and for every u in $\text{Fut}_{\mathcal{A}}(s)$ there exists a unique v in $\text{Fut}_{\mathcal{B}'}(s)$ such that $u \in \varphi(v)$.*

Proof. Suppose by way of contradiction that there exist v and v' in $\text{Fut}_{\mathcal{B}'}(s)$ such that $u \in \varphi(v) \cap \varphi(v')$. We write $v = x_n \dots x_2 x_1$ and $v' = x'_m \dots x'_2 x'_1$. There exist $w_i \in \varphi(x_i)$ and $w'_i \in \varphi(x'_i)$ such that $u = w_n \dots w_2 w_1 = w'_m \dots w'_2 w'_1$.

Let j be the smallest index such that either $w_j \neq w'_j$ or $x_j \neq x'_j$. Let r be the unique state of \mathcal{B}' such that $x_{j-1} \dots x_1 \in \text{Fut}_{\mathcal{B}'}(r)$. By Property 3.2, $\varphi(J_r)$ is a suffix code and then $w_j = w'_j$, and $\varphi(x_j) \cap \varphi(x'_j) \neq \emptyset$ implies $x_j = x'_j$. \square

Let K be a derived term in $D(E)$. As \mathcal{A} is co-deterministic, and by Proposition 2.8 either there exists a unique state $s \neq q$ in \mathcal{A} such that $|K| \subseteq \text{Fut}_{\mathcal{A}}(s)$ — **Case 1** — or $|K| \subseteq \text{Fut}_{\mathcal{A}}(q)$ — **Case 2**.

Lemma 3.4. *Let K in $D(E)$. There exists K' in $D(E')$ such that either:*

Case 1 *there exists $s \neq q$ such that $|K'| \subseteq \text{Fut}_{\mathcal{B}'}(s)$ and $K = \varphi(K')$, or*

Case 2 *there exists a successor r of q such that $|K'| \subseteq \text{Fut}_{\mathcal{B}'}(r)$ and $K = G^*H_r\varphi(K')$ or $K = H_r\varphi(K')$ (if G is null).*

Proof. From (11), K is either (**Case 1**) in $\varphi(D(E'))$ or (**Case 2**) in

$$\bigcup_{c_{p,r} \in C} \left[\bigcup_{a \in F_p} \frac{\partial}{\partial a} \varphi(c_{p,r}) \right] \varphi(D_{c_{p,r}}(E')) .$$

In **Case 1**, there exists K' such that $K = \varphi(K')$, by Proposition 2.8 there exists s in \mathcal{B}' such that $|K'| \subseteq \text{Fut}_{\mathcal{B}'}(s)$ and by Property 3.1 $|K| \subseteq \text{Fut}_{\mathcal{A}}(s)$.

In **Case 2**, there exist $c_{p,r} \in C$, $a \in A$ and $K' \in D_{c_{p,r}}(E')$ such that

$$K = \left[\frac{\partial}{\partial a} \varphi(c_{p,r}) \right] \varphi(K') = G^*H_r\varphi(K') .$$

As $K' \in D_{c_{p,r}}(E')$, $|K'| \subseteq \text{Fut}_{\mathcal{B}'}(r)$. By Property 3.1 $\varphi(|K'|) \subseteq \text{Fut}_{\mathcal{A}}(r)$ and thus, by definition of G and H_r , $|K| \subseteq \text{Fut}_{\mathcal{A}}(q)$. \square

Proposition 3.5. *Let \mathcal{A} be a normalized co-deterministic automaton, and $E = \Phi(\mathcal{A})$. Then, the derived term automaton $\Delta(E)$ of E is co-deterministic.*

Proof. A normalized automaton has at least two states. If \mathcal{A} has only two states, $\Phi(\mathcal{A})$ is reduced to a sum of letters and $\Delta(\Phi(\mathcal{A}))$ is clearly co-deterministic.

By induction, $\Delta(E')$ is co-deterministic, which means, by Property 1.5, that the interpretations of the derived terms of E' are pairwise disjoint. We prove now that the interpretations of elements of $D(E)$ are disjoint, which implies that $\Delta(E)$ is co-deterministic.

Let K_1 and K_2 be two distinct derived terms of $D(E)$ and assume that there exists u in $|K_1| \cap |K_2|$. By (11) and Lemma 3.4, there exist K'_1 and K'_2 in $D(E')$ such that one of the following three cases holds:

Case 1 $K_1 = \varphi(K'_1)$ and $K_2 = \varphi(K'_2)$.

Case 2.1 $K_1 = \varphi(K'_1)$ and $K_2 = G^*H_r\varphi(K'_2)$.

Case 2.2 $K_1 = G^*H_{r_1}\varphi(K'_1)$ and $K_2 = G^*H_{r_2}\varphi(K'_2)$.

Case 1 From Lemma 3.4 there exist states s_1 and s_2 such that $|K'_1| \subseteq \text{Fut}_{\mathcal{B}'}(s_1)$ and $|K'_2| \subseteq \text{Fut}_{\mathcal{B}'}(s_2)$, and $|K_1| \subseteq \text{Fut}_{\mathcal{A}}(s_1)$ and $|K_2| \subseteq \text{Fut}_{\mathcal{A}}(s_2)$. As \mathcal{A} is co-deterministic, $s_1 = s_2 = s \neq q$ and by Lemma 3.3, there exists a unique v in $\text{Fut}_{\mathcal{B}'}(s)$ such that $u \in \varphi(v)$. Hence $v \in |K'_1| \cap |K'_2|$. The induction hypothesis implies $K'_1 = K'_2$ and thus $K_1 = K_2$.

Case 2.1 From Lemma 3.4 there exist states $s_1 \neq q$ such that $|K_1| \subseteq \text{Fut}_{\mathcal{A}}(s_1)$ and $|K_2| \subseteq \text{Fut}_{\mathcal{A}}(q)$. As \mathcal{A} is co-deterministic, $|K_1| \cap |K_2| = \emptyset$.

Case 2.2 From Lemma 3.4, there exist successors of q , r_1 and r_2 such that $|K'_1| \subseteq \text{Fut}_{\mathcal{B}'}(r_1)$ and $|K'_2| \subseteq \text{Fut}_{\mathcal{B}'}(r_2)$, and $|K_1| \subseteq \text{Fut}_{\mathcal{A}}(q)$ and $|K_2| \subseteq \text{Fut}_{\mathcal{A}}(q)$. Hence $u = g_1h_1w_1 = g_2h_2w_2$, with $g_1, g_2 \in G^*$, $h_1 \in H_{r_1}$, $h_2 \in H_{r_2}$, $w_1 \in \varphi(|K'_1|)$ and $w_2 \in \varphi(|K'_2|)$. As $G \cap H_{r_1} = G \cap H_{r_2} = \emptyset$, $g_1 = g_2$, $h_1 = h_2$, and thus $w_1 = w_2$. Since \mathcal{A} is co-deterministic, there is a unique state s , different from q , such that w is in $\text{Fut}_{\mathcal{A}}(s)$. Hence, by Case 1, $\varphi(K'_1) = \varphi(K'_2)$ and thus $K_1 = K_2$. \square

Corollary 3.6. *Let \mathcal{A} be a normalized co-deterministic automaton, and $E = \Phi(\mathcal{A})$. Then, the broken derived term automaton $\Delta'(E)$ of E is co-deterministic.*

Proof. As above, the corollary trivially holds if \mathcal{A} has only two states.

By induction, $\Delta'(E')$ is co-deterministic, which means, by Property 1.5, that the interpretations of the broken derived terms of E' are pairwise disjoint. We prove now that the interpretations of elements of $\text{BD}(E)$ are disjoint, which implies that $\Delta'(E)$ is co-deterministic.

By Equation (5), $\text{BD}(E) = d(D(E))$ and, by Proposition 3.5, the interpretations of the derived terms of E are disjoint. If there exists a word u in the intersection of the interpretation of two broken derived terms K_1 and K_2 , then, there exists a *derived term* L of E , such that both K_1 and K_2 are in $d(L)$.

Case 1 There exists L' in $D(E')$ such that $L = \varphi(L')$. Then, there exists K'_1 and K'_2 in $d(L')$ such that $K_1 = \varphi(K'_1)$ and $K_2 = \varphi(K'_2)$. By Lemma 3.4 and Lemma 3.3, there exists a unique word v in $|L'|$ such that $\varphi(v) = u$, hence, v is in $|K'_1| \cap |K'_2|$. By induction hypothesis, broken derived terms of E' have disjoint interpretations, thus $K'_1 = K'_2$ and $K_1 = K_2$.

Case 2 There exist L' in $D(E')$, G and H_r , such that $L = G^*H_r\varphi(L')$ or $L = H_r\varphi(L')$ (if G is null).

In the first subcase, $d(L) = \{L\}$, hence $K_1 = K_2 = L$; in the second subcase, $K_1 = a_1\varphi(L')$ and $K_2 = a_2\varphi(L')$, a_1 and a_2 are both the first letter of u and are therefore equal, hence $K_1 = K_2$. \square

Theorem 1.6 is now established. Indeed, either \mathcal{A} is normalized, and nothing is to be added, or we apply Corollary 3.6 to \mathcal{A}_\S and as the interpretations of any two broken derived terms of E_\S are disjoint so are the corresponding broken derived terms of E : the broken derived term automaton is co-deterministic. \square

CONCLUSION

As we wanted to be complete, and not leave another flaw behind us, the proof of the theorem is long, much longer than we expected, and somewhat laboured. It goes by a repeated forth and back between the properties of the terms and of their interpretation. It is also based on an interplay between the derived terms and the broken ones, which let us think that the truth may be hidden somewhere in between. . .

In any case, this proof brought under light these broken terms that we defined in [6] in an exploratory way. They deserve to be further studied, a task that we have already begun [1].

Acknowledgements The authors are pleased to thanks Florent Terrones, from the Vaucanson Group, who pointed out an example which led to this correction and to Pierre-Yves Angrand who drew our attention on a problem with the breaking operation $d()$ and gave the correct definition.

REFERENCES

- [1] P.-Y. ANGRAND, S. LOMBARDY AND J. SAKAROVITCH, On the broken derived terms of a rational expression. *In preparation*.
- [2] V. ANTIMIROV, Partial derivatives of regular expressions and finite automaton constructions. *Theoret. Computer Sci.* **155** (1996), 291–319.
- [3] J. A. BRZozowski, Derivatives of regular expressions. *J. Assoc. Comput. Mach.* **11** (1964), 481–494.
- [4] P. CARON AND M. FLOURET, Glushkov construction for series: the non commutative case, *Int. J. Comput. Math.* **80** (2003), 457–472.
- [5] V. GLUSHKOV, The abstract theory of automata. *Russian Mathematical Surveys* **16** (1961), 1–53.
- [6] S. LOMBARDY AND J. SAKAROVITCH, Derivatives of rational expressions with multiplicity, *Theoret. Computer Sci.* **332** (2005), 141–177. (Journal version of *Proc. MFCS 02*, LNCS 2420 (2002), 471–482.)
- [7] S. LOMBARDY AND J. SAKAROVITCH, How expressions can code for automata, *RAIRO – Theoret. Informatics and Applications* **39** (2005), 217–237. (Journal version of *Proc. of LATIN 2004*, LNCS 2976 (2004), 242–251.)