

Derived-Term Automata of Multitape Rational Expressions

Akim Demaille `akim@lrde.epita.fr`

EPITA Research and Development Laboratory (LRDE)
14-16, rue Voltaire, 94276 Le Kremlin-Bicêtre, France

Abstract. We consider (weighted) rational expressions to denote series over Cartesian products of monoids. We define an operator $|$ to build multitape expressions such as $(a^+ | x + b^+ | y)^*$. We introduce *expansions*, which generalize the concept of derivative of a rational expression, but relieved from the need of a free monoid. We propose an algorithm based on expansions to build multitape automata from multitape expressions.

1 Introduction

Automata and rational (or regular) expressions share the same expressive power, with algorithms going from one to the other. This fact made rational expressions an extremely handy practical tool to specify some rational languages in a concise way, from which acceptors (automata) are built. There are many largely used implementations, probably starting with Ken Thompson [15], the creator of Unix, `grep`, etc.

There are numerous algorithms to build an automaton from an expression. We are particularly interested in the derivative-based family of algorithms [3–5, 7, 10], because they offer a very natural interpretation to states (they are labeled by an expression that denotes the future of the states, i.e., the language/series accepted from this state). This allowed to support several extensions: extended operators (intersection, complement) [4, 5], weights [10], additional products (shuffle, infiltration), etc.

Multitape automata, including transducers, share many properties with “single-tape” automata, in particular the Fundamental Theorem [14, Theorem 2.1, p. 409]: under appropriate conditions, multitape automata and rational (multitape) series share the same expressive power. However, as far as the author knows, there is no definition of multitape rational expressions that allows expressions such as $E_2 := (a^+ | x + b^+ | y)^*$ (Example 5). To denote such a binary relation between words, one had to build a (usual) rational expression in “normal form”, without tupling of expressions but only tuples of letters such as a set of generators. So for instance instead of E_2 , one must use $E'_2 := ((a | \varepsilon)^+(\varepsilon | x) + (b | \varepsilon)^+(\varepsilon | y))^*$, which is larger, as is its derived-term automaton.

The contributions of this paper are twofold: we define (weighted) multitape rational expressions featuring a $|$ operator, and we provide an algorithm to build an equivalent automaton. This algorithm is a generalization of the derived-term based algorithms, freed from the requirement that the monoid is free.

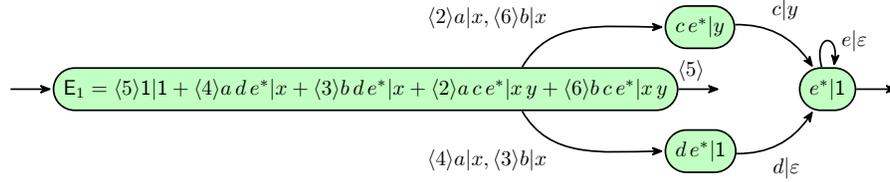


Fig. 1. The derived-term automaton of E_1 (see Examples 1 to 3) with $E_1 := \langle 5 \rangle 1|1 + \langle 4 \rangle a d e^*|x + \langle 3 \rangle b d e^*|x + \langle 2 \rangle a c e^*|x y + \langle 6 \rangle b c e^*|x y$.

We first settle the notations in Sect. 2, provide an algorithm to compute the expansion of an expression in Sect. 3, which is used in Sect. 4 to propose an alternative construction of the derived-term automaton.

The constructs exposed in this paper are implemented in Vcsn¹. Vcsn is a free-software platform dedicated to weighted automata and rational expressions [8]; its lowest layer is a C++ library, on top of which Python/IPython bindings provide an interactive graphical environment.

2 Notations

Our purpose is to define (weighted) multitape rational expressions, such as $E_1 := \langle 5 \rangle 1|1 + \langle 4 \rangle a d e^*|x + \langle 3 \rangle b d e^*|x + \langle 2 \rangle a c e^*|x y + \langle 6 \rangle b c e^*|x y$ (weights are written in angle brackets). It relates ade with x , with weight 4. We introduce an algorithm to build a multitape automaton (aka *transducer*) from such an expression, e.g., Fig. 1. This algorithm relies on *rational expansions*. They are to the derivatives of rational expressions what differential forms are to the derivatives of functions. Defining expansions requires several concepts, defined bottom-up in this section. The following figure presents these different entities, how they relate to each other, and where we are heading to: given a weighted multitape rational expression such as E_1 , compute *its* expansion:

$$\underbrace{\underbrace{\langle 5 \rangle}_{\text{Weight}} \oplus \underbrace{a|x}_{\text{Label}} \odot \left[\underbrace{\langle 2 \rangle \odot ce^*|y}_{\text{Derived term}} \oplus \underbrace{\langle 4 \rangle \odot de^*|1}_{\text{Monomial}} \right]}_{\text{First}} \oplus b|x \odot \underbrace{\left[\langle 6 \rangle \odot ce^*|y \oplus \langle 3 \rangle \odot de^*|1 \right]}_{\text{Polynomial (Sect. 2.3)}}}_{\text{Proper part of the expansion}}$$

Expansion (Sect. 2.4)

from which we build its derived-term automaton (Fig. 1).

It is helpful to think of expansions as a normal form for expressions.

¹ See the interactive environment, <http://vcsn-sandbox.lrde.epita.fr>, or its documentation, http://vcsn.lrde.epita.fr/dload/2.3/notebooks/expression.derived_term.html, or this paper's companion notebook, <http://vcsn.lrde.epita.fr/dload/2.3/notebooks/CIAA-2016.html>.

2.1 Rational Series

Series will be used to define the semantics of the forthcoming structures: they are to weighted automata what languages are to Boolean automata. Not all languages are rational (denoted by an expression), and similarly, not all series are rational (denoted by a weighted expression). We follow Sakarovitch [14, Chap. III].

In order to cope with (possibly) several tapes, we cannot rely on the traditional definitions based on the free monoid A^* for some alphabet A .

Labels Let M be a monoid (e.g., A^* or $A^* \times B^*$), whose neutral element is denoted ε_M , or ε when clear from the context. For consistency with the way transducers are usually represented, we use $m \mid n$ rather than (m, n) to denote the pair of m and n . For instance $\varepsilon_{A^* \times B^*} = \varepsilon_{A^*} \mid \varepsilon_{B^*}$, and $\varepsilon_M \mid a \in M \times \{a\}^*$. A *set of generators* G of M is a subset of M such that $G^* = M$. A monoid M is of *finite type* (or *finitely generated*) if it admits a finite set of generators. A monoid M is *graded* if it admits a *gradation* function $|\cdot| \in M \rightarrow \mathbb{N}$ such that $\forall m, n \in M$, $|m| = 0$ iff $m = \varepsilon$, and $|mn| = |m| + |n|$. Cartesian products of graded monoids are graded, and Cartesian products of finitely generated monoids are finitely generated. Free monoids and Cartesian products of free monoids are graded and finitely generated.

Weights Let $\langle \mathbb{K}, +, \cdot, 0_{\mathbb{K}}, 1_{\mathbb{K}} \rangle$ (or \mathbb{K} for short) be a semiring whose (possibly non commutative) multiplication will be denoted by juxtaposition. \mathbb{K} is *commutative* if its multiplication is. \mathbb{K} is a *topological semiring* if it is equipped with a topology, and both addition and multiplication are continuous. It is *strong* if the product of two summable families is summable.

Series A (formal power) *series* over M with *weights* (or *multiplicities*) in \mathbb{K} is a map from M to \mathbb{K} . The weight of $m \in M$ in a series s is denoted $s(m)$. The *null series*, $m \mapsto 0_{\mathbb{K}}$, is denoted 0 ; for any $m \in M$ (including ε_M), m denotes the series $u \mapsto 1_{\mathbb{K}}$ if $u = m$, $0_{\mathbb{K}}$ otherwise. If M is of finite type, then we can define the Cauchy product of series. $s \cdot t := m \mapsto \sum_{u, v \in M \mid uv=m} s(u) \cdot t(v)$. Equipped with the pointwise addition ($s + t := m \mapsto s(m) + t(m)$) and \cdot as multiplication, the set of these series forms a semiring denoted $\langle \mathbb{K}\langle\langle M \rangle\rangle, +, \cdot, 0, \varepsilon \rangle$.

The *constant term* of a series s , denoted s_ε , is $s(\varepsilon)$, the weight of the empty word. A series s is *proper* if $s_\varepsilon = 0_{\mathbb{K}}$. The *proper part* of s is the proper series s_p such that $s = s_\varepsilon + s_p$.

Star The *star* of a series is an infinite sum: $s^* := \sum_{n \in \mathbb{N}} s^n$. To ensure semantic soundness, we need M to be graded monoid and \mathbb{K} to be a strong topological semiring.

Proposition 1. *Let M be a graded monoid and \mathbb{K} a strong topological semiring. Let $s \in \mathbb{K}\langle\langle M \rangle\rangle$, s^* is defined iff s_ε^* is defined and then $s^* = s_\varepsilon^* + s_\varepsilon^* s_p s^*$.*

Proof. By [14, Prop. 2.6, p. 396] s^* is defined iff s_ε^* is defined and then $s^* = (s_\varepsilon^* s_p)^* s_\varepsilon^* = s_\varepsilon^* (s_p s_\varepsilon^*)^*$. The result then follows directly from $s^* = \varepsilon + s s^*$: $s^* = s_\varepsilon^* (s_p s_\varepsilon^*)^* = s_\varepsilon^* (\varepsilon + (s_p s_\varepsilon^*) (s_p s_\varepsilon^*)^*) = s_\varepsilon^* + s_\varepsilon^* s_p (s_\varepsilon^* (s_p s_\varepsilon^*)^*) = s_\varepsilon^* + s_\varepsilon^* s_p s^*$. \square

Tuple We suppose \mathbb{K} is commutative. The *tupling* of two series $s \in \mathbb{K}\langle\langle M \rangle\rangle, t \in \mathbb{K}\langle\langle N \rangle\rangle$, is the series $s \mid t := m \mid n \in M \times N \mapsto s(m)t(n)$. It is a member of $\mathbb{K}\langle\langle M \times N \rangle\rangle$.

Proposition 2. For all series $s, s' \in \mathbb{K}\langle\langle M \rangle\rangle$ and $t, t' \in \mathbb{K}\langle\langle N \rangle\rangle$, $(s + s') | t = s | t + s' | t$ and $s | (t + t') = s | t + s | t'$.

Proof. Let $m | n \in M \times N$. $((s + s') | t)(m | n) = (s + s')(m) \cdot t(n) = (s(m) + s'(m)) \cdot t(n) = s(m) \cdot t(n) + s'(m) \cdot t(n) = (s | t)(m | n) + (s' | t)(m | n) = (s | t + s' | t)(m | n)$. Likewise for right distributivity. \square

From now on, M is a graded monoid of finite type, and \mathbb{K} a commutative strong topological semiring.

2.2 Weighted Rational Expressions

Contrary to the usual definition, we do not require a finite alphabet: any set of generators $G \subseteq M$ will do. For expressions with more than one tape, we required \mathbb{K} to be commutative; however, for single tape expressions, our results apply to non-commutative semirings, hence there are two exterior products.

Definition 1 (Expression). A rational expression E over G is a term built from the following grammar, where $a \in G$ denotes any non empty label, and $k \in \mathbb{K}$ any weight: $E ::= 0 \mid 1 \mid a \mid E + E \mid \langle k \rangle E \mid E \langle k \rangle \mid E \cdot E \mid E^* \mid E | E$.

Expressions are syntactic; they are finite notations for (some) series.

Definition 2 (Series Denoted by an Expression). Let E be an expression. The series denoted by E , noted $\llbracket E \rrbracket$, is defined by induction on E :

$$\begin{aligned} \llbracket 0 \rrbracket &:= 0 & \llbracket 1 \rrbracket &:= \varepsilon & \llbracket a \rrbracket &:= a & \llbracket E + F \rrbracket &:= \llbracket E \rrbracket + \llbracket F \rrbracket & \llbracket \langle k \rangle E \rrbracket &:= k \llbracket E \rrbracket \\ \llbracket E \langle k \rangle \rrbracket &:= \llbracket E \rrbracket k & \llbracket E \cdot F \rrbracket &:= \llbracket E \rrbracket \cdot \llbracket F \rrbracket & \llbracket E^* \rrbracket &:= \llbracket E \rrbracket^* & \llbracket E | F \rrbracket &:= \llbracket E \rrbracket | \llbracket F \rrbracket \end{aligned}$$

An expression is *valid* if it denotes a series. More specifically, there are two requirements. First, the expression must be well-formed, i.e., concatenation and disjunction must be applied to expressions of appropriate number of tapes. For instance, $a + b | c$ and $a(b | c)$ are ill-formed, $(a | b)^* | c + a | (b | c)^*$ is well-formed. Second, to ensure that $\llbracket F \rrbracket^*$ is well defined for each subexpression of the form F^* , the constant term of $\llbracket F \rrbracket$ must be *starrable* in \mathbb{K} (Proposition 1). This definition, which involves series (semantics) to define a property of expressions (syntax), will be made effective (syntactic) with the appropriate definition of the constant term $d_\varepsilon(E)$ of an expression E (Definition 6).

Let $[n]$ denote $\{1, \dots, n\}$. The *size* (aka *length*) of a (valid) expression E , $|E|$, is its total number of symbols, not counting parenthesis; for a given tape number $i \in [k]$ the *width on tape i* , $\|E\|_i$, is the number of occurrences of labels on the tape i , the *width* of E (aka *literal length*), $\|E\| := \sum_{i \in [k]} \|E\|_i$ is the total number of occurrences of labels.

Two expressions E and F are *equivalent* iff $\llbracket E \rrbracket = \llbracket F \rrbracket$. Some expressions are “trivially equivalent”; any candidate expression will be rewritten via the following *trivial identities*. Any subexpression of a form listed to the left of a ‘ \Rightarrow ’ is rewritten as indicated on the right.

$$E + 0 \Rightarrow E \quad 0 + E \Rightarrow E$$

$$\begin{aligned}
 \langle 0_{\mathbb{K}} \rangle E &\Rightarrow 0 & \langle 1_{\mathbb{K}} \rangle E &\Rightarrow E & \langle k \rangle 0 &\Rightarrow 0 & \langle k \rangle \langle h \rangle E &\Rightarrow \langle kh \rangle E \\
 E \langle 0_{\mathbb{K}} \rangle &\Rightarrow 0 & E \langle 1_{\mathbb{K}} \rangle &\Rightarrow E & 0 \langle k \rangle &\Rightarrow 0 & E \langle k \rangle \langle h \rangle &\Rightarrow E \langle kh \rangle \\
 \langle \langle k \rangle E \rangle \langle h \rangle &\Rightarrow \langle k \rangle \langle E \rangle \langle h \rangle & \ell \langle k \rangle &\Rightarrow \langle k \rangle \ell \\
 E \cdot 0 &\Rightarrow 0 & 0 \cdot E &\Rightarrow 0 \\
 \langle \langle k \rangle^? 1 \rangle \cdot E &\Rightarrow \langle k \rangle^? E & E \cdot \langle \langle k \rangle^? 1 \rangle &\Rightarrow E \langle k \rangle^? \\
 0^* &\Rightarrow 1 \\
 \langle \langle k \rangle^? E \rangle \mid \langle \langle h \rangle^? F \rangle &\Rightarrow \langle kh \rangle^? E \mid F
 \end{aligned}$$

where E is a rational expression, $\ell \in G \cup \{1\}$ a label, $k, h \in \mathbb{K}$ weights, and $\langle k \rangle^? \ell$ denotes either $\langle k \rangle \ell$, or ℓ in which case $k = 1_{\mathbb{K}}$ in the right-hand side of \Rightarrow . The choice of these identities is beyond the scope of this paper (see [14]), however note that they are limited to trivial properties; in particular *linearity* (“weighted ACI”: associativity, commutativity and $\langle k \rangle E + \langle h \rangle E \Rightarrow \langle k + h \rangle E$) is not enforced. In practice, additional identities help reducing the automaton size [12].

2.3 Rational Polynomials

At the core of the idea of “partial derivatives” introduced by Antimirov [3], is that of *sets* of rational expressions, later generalized in *weighted sets* by Lombardy and Sakarovitch [10], i.e., functions (partial, with finite domain) from the set of rational expressions into $\mathbb{K} \setminus \{0_{\mathbb{K}}\}$. It proves useful to view such structures as “polynomials of expressions”. In essence, they capture the linearity of addition.

Definition 3 (Rational Polynomial). *A polynomial (of rational expressions) is a finite (left) linear combination of expressions. Syntactically it is a term built from the grammar $P ::= 0 \mid \langle k_1 \rangle \odot E_1 \oplus \dots \oplus \langle k_n \rangle \odot E_n$ where $k_i \in \mathbb{K} \setminus \{0_{\mathbb{K}}\}$ denote non-null weights, and E_i denote non-null expressions. Expressions may not appear more than once in a polynomial. A monomial is a pair $\langle k_i \rangle \odot E_i$.*

We use specific symbols (\odot and \oplus) to clearly separate the outer polynomial layer from the inner expression layer. Let $P = \bigoplus_{i \in [n]} \langle k_i \rangle \odot E_i$ be a polynomial of expressions. The “*projection*” of P is the expression $\text{expr}(P) := \langle k_1 \rangle E_1 + \dots + \langle k_n \rangle E_n$ (or 0 if P is null); this operation is performed on a canonical form of the polynomial (expressions are sorted in a well defined order). Polynomials denote series: $\llbracket P \rrbracket := \llbracket \text{expr}(P) \rrbracket$. The *terms* of P is the set $\text{exprs}(P) := \{E_1, \dots, E_n\}$.

Example 1. Let $E_1 := \langle 5 \rangle 1 \mid 1 + \langle 4 \rangle a d e^* \mid x + \langle 3 \rangle b d e^* \mid x + \langle 2 \rangle a c e^* \mid x y + \langle 6 \rangle b c e^* \mid x y$. Polynomial ‘ $P_{1,a \mid x} := \langle 2 \rangle \odot c e^* \mid y \oplus \langle 4 \rangle \odot d e^* \mid 1$ ’ has two monomials: ‘ $\langle 2 \rangle \odot c e^* \mid y$ ’ and ‘ $\langle 4 \rangle \odot d e^* \mid 1$ ’. It denotes the (left) quotient of $\llbracket E_1 \rrbracket$ by $a \mid x$, and ‘ $P_{1,b \mid x} := \langle 6 \rangle \odot c e^* \mid y \oplus \langle 3 \rangle \odot d e^* \mid 1$ ’ the quotient by $b \mid x$.

Let $P = \bigoplus_{i \in [n]} \langle k_i \rangle \odot E_i$, $Q = \bigoplus_{j \in [m]} \langle h_j \rangle \odot F_j$ be polynomials, k a weight and F an expression, all possibly null, we introduce the following operations:

$$P \cdot F := \bigoplus_{i \in [n]} \langle k_i \rangle \odot (E_i \cdot F) \quad \langle k \rangle P := \bigoplus_{i \in [n]} \langle k k_i \rangle \odot E_i \quad P \langle k \rangle := \bigoplus_{i \in [n]} \langle k_i \rangle \odot (E_i \langle k \rangle)$$

$$\begin{aligned} P \mid 1 &:= \bigoplus_{i \in [n]} \langle k_i \rangle \odot E_i \mid 1 & 1 \mid P &:= \bigoplus_{i \in [n]} \langle k_i \rangle \odot 1 \mid E_i \\ P \mid Q &:= \bigoplus_{(i,j) \in [n] \times [m]} \langle k_i \cdot h_j \rangle \odot E_i \mid F_j \end{aligned}$$

Trivial identities might simplify the result. Note the asymmetry between left and right exterior products. The addition of polynomials is commutative, multiplication by zero (be it an expression or a weight) evaluates to the null polynomial, and the left-multiplication by a weight is distributive.

Lemma 1. $\llbracket P \cdot F \rrbracket = \llbracket P \rrbracket \cdot \llbracket F \rrbracket$ $\llbracket \langle k \rangle P \rrbracket = \langle k \rangle \llbracket P \rrbracket$ $\llbracket P \langle k \rangle \rrbracket = \llbracket P \rrbracket \langle k \rangle$
 $\llbracket P \mid Q \rrbracket = \llbracket P \rrbracket \mid \llbracket Q \rrbracket$.

2.4 Rational Expansions

Definition 4 (Rational Expansion). A rational expansion X is a term $X := \langle X_\varepsilon \rangle \oplus a_1 \odot [X_{a_1}] \oplus \cdots \oplus a_n \odot [X_{a_n}]$ where $X_\varepsilon \in \mathbb{K}$ is a weight (possibly null), $a_i \in G \setminus \{\varepsilon\}$ non-empty labels (occurring at most once), and X_{a_i} non-null polynomials. The constant term is X_ε , the proper part is $X_p := a_1 \odot [X_{a_1}] \oplus \cdots \oplus a_n \odot [X_{a_n}]$, the firsts is $f(X) := \{a_1, \dots, a_n\}$ (possibly empty) and the terms $\text{exprs}(X) := \bigcup_{i \in [n]} \text{exprs}(X_{a_i})$.

To ease reading, polynomials are written in square brackets. Contrary to expressions and polynomials, there is no specific term for the null expansion: it is represented by $\langle 0_{\mathbb{K}} \rangle$, the null weight. Except for this case, null constant terms are left implicit. Expansions will be written: $X = \langle X_\varepsilon \rangle \oplus \bigoplus_{a \in f(X)} a \odot [X_a]$. When more convenient, we write $X(\ell)$ instead of X_ℓ for $\ell \in f(X) \cup \{\varepsilon\}$.

An expansion X can be “projected” as a rational expression $\text{expr}(X)$ by mapping weights, labels and polynomials to their corresponding rational expressions, and \oplus/\odot to the sum/concatenation of expressions. Again, this is performed on a canonical form of the expansion: labels are sorted. Expansions also denote series: $\llbracket X \rrbracket := \llbracket \text{expr}(X) \rrbracket$. An expansion X is *equivalent* to an expression E iff $\llbracket X \rrbracket = \llbracket E \rrbracket$.

Example 2 (Example 1 continued). Expansion $X_1 := \langle 5 \rangle \oplus a|x \odot [P_{1,a|x}] \oplus b|x \odot [P_{1,b|x}]$ has $X_1(\varepsilon) = \langle 5 \rangle$ as constant term, and maps the generator $a|x$ (resp. $b|x$) to the polynomial $X_1(a|x) = P_{1,a|x}$ (resp. $X_1(b|x) = P_{1,b|x}$). X_1 can be proved to be equivalent to E_1 .

Let X, Y be expansions, k a weight, and E an expression (all possibly null):

$$X \oplus Y := \langle X_\varepsilon + Y_\varepsilon \rangle \oplus \bigoplus_{a \in f(X) \cup f(Y)} a \odot [X_a \oplus Y_a] \quad (1)$$

$$\langle k \rangle X := \langle k X_\varepsilon \rangle \oplus \bigoplus_{a \in f(X)} a \odot [\langle k \rangle X_a] \quad X \langle k \rangle := \langle X_\varepsilon k \rangle \oplus \bigoplus_{a \in f(X)} a \odot [X_a \langle k \rangle] \quad (2)$$

$$X \cdot E := \bigoplus_{a \in f(X)} a \odot [X_a \cdot E] \quad \text{with } X \text{ proper: } X_\varepsilon = 0_{\mathbb{K}} \quad (3)$$

$$\begin{aligned}
 X | Y := & \langle X_\varepsilon Y_\varepsilon \rangle \oplus \langle X_\varepsilon \rangle \bigoplus_{b \in f(Y)} (\varepsilon | b) \odot (1 | Y_b) \oplus \langle Y_\varepsilon \rangle \bigoplus_{a \in f(X)} (a | \varepsilon) \odot (X_a | 1) \\
 & \oplus \bigoplus_{a | b \in f(X) \times f(Y)} (a | b) \odot (X_a | Y_b)
 \end{aligned} \tag{4}$$

Since by definition expansions never map to null polynomials, some firsts might be smaller than suggested by these equations. For instance in \mathbb{Z} the sum of $\langle 1 \rangle \oplus a \odot [\langle 1 \rangle \odot b]$ and $\langle 1 \rangle \oplus a \odot [\langle -1 \rangle \odot b]$ is $\langle 2 \rangle$.

The following lemma is simple to establish: lift semantic equivalences, such as [Proposition 2](#), to syntax, using [Lemma 1](#).

$$\begin{aligned}
 \llbracket X \oplus Y \rrbracket &= \llbracket X \rrbracket + \llbracket Y \rrbracket & \llbracket \langle k \rangle X \rrbracket &= \langle k \rangle \llbracket X \rrbracket & \llbracket X \langle k \rangle \rrbracket &= \llbracket X \rrbracket \langle k \rangle \\
 \llbracket X \cdot E \rrbracket &= \llbracket X \rrbracket \cdot \llbracket E \rrbracket & \llbracket X | Y \rrbracket &= \llbracket X \rrbracket | \llbracket Y \rrbracket
 \end{aligned}$$

2.5 Finite Weighted Automata

Definition 5 (Weighted Automaton). A weighted automaton \mathcal{A} is a tuple $\langle M, G, \mathbb{K}, Q, E, I, T \rangle$ where:

- M is a monoid,
- G (the labels) is a set of generators of M ,
- \mathbb{K} (the set of weights) is a semiring,
- Q is a finite set of states,
- I and T are the initial and final functions from Q into \mathbb{K} ,
- E is a (partial) function from $Q \times G \times Q$ into $\mathbb{K} \setminus \{0_{\mathbb{K}}\}$;
its domain represents the transitions: (source, label, destination).

An automaton is proper if no label is ε_M .

A computation $p = (q_0, a_0, q_1)(q_1, a_1, q_2) \cdots (q_n, a_n, q_{n+1})$ in an automaton is a sequence of transitions where the source of each is the destination of the previous one; its label is $a_0 a_1 \cdots a_n \in M$, its weight is $I(q_0) \otimes E(q_0, a_0, q_1) \otimes \cdots \otimes E(q_n, a_n, q_{n+1}) \otimes T(q_{n+1}) \in \mathbb{K}$. The evaluation of word u by \mathcal{A} , $\mathcal{A}(u)$, is the sum of the weights of all the computations labeled by u , or $0_{\mathbb{K}}$ if there are none. The behavior of an automaton \mathcal{A} is the series $\llbracket \mathcal{A} \rrbracket := m \mapsto \mathcal{A}(m)$. A state q is *initial* if $I(q) \neq 0_{\mathbb{K}}$. A state q is *accessible* if there is a computation from an initial state to q . The *accessible* part of an automaton \mathcal{A} is the subautomaton whose states are the accessible states of \mathcal{A} . The size of an automaton, $|\mathcal{A}|$, is its number of states.

We are interested, given an expression E , in an algorithm to compute an automaton \mathcal{A}_E such that $\llbracket \mathcal{A}_E \rrbracket = \llbracket E \rrbracket$ ([Definition 7](#)). To this end, we first introduce a simple recursive procedure to compute the expansion of an expression.

3 Expansion of a Rational Expression

Definition 6 (Expansion of a Rational Expression). The expansion of a rational expression E , written $d(E)$, is defined inductively as follows:

$$d(0) := \langle 0_{\mathbb{K}} \rangle \quad d(1) := \langle 1_{\mathbb{K}} \rangle \quad d(a) := a \odot [\langle 1_{\mathbb{K}} \rangle \odot 1] \tag{5}$$

$$d(\mathbf{E} + \mathbf{F}) := d(\mathbf{E}) \oplus d(\mathbf{F}) \quad (6)$$

$$d(\langle k \rangle \mathbf{E}) := \langle k \rangle d(\mathbf{E}) \quad d(\mathbf{E} \langle k \rangle) := d(\mathbf{E}) \langle k \rangle \quad (7)$$

$$d(\mathbf{E} \cdot \mathbf{F}) := d_p(\mathbf{E}) \cdot \mathbf{F} \oplus \langle d_\varepsilon(\mathbf{E}) \rangle d(\mathbf{F}) \quad (8)$$

$$d(\mathbf{E}^*) := \langle d_\varepsilon(\mathbf{E}^*) \rangle \oplus \langle d_\varepsilon(\mathbf{E}^*) \rangle d_p(\mathbf{E}) \cdot \mathbf{E}^* \quad (9)$$

$$d(\mathbf{E} \mid \mathbf{F}) := d(\mathbf{E}) \mid d(\mathbf{F}) \quad (10)$$

where $d_\varepsilon(\mathbf{E}) := d(\mathbf{E})_\varepsilon$, $d_p(\mathbf{E}) := d(\mathbf{E})_p$ are the constant term/proper part of $d(\mathbf{E})$.

The right-hand sides are indeed expansions. The computation trivially terminates: induction is performed on strictly smaller subexpressions. These formulas are enough to compute the expansion of an expression; there is no secondary process to compute the firsts — indeed $d(a) := a \odot [\langle 1_{\mathbb{K}} \rangle \odot 1]$ suffices and every other case simply propagates or assembles the firsts — or the constant terms. Note that the firsts are a subset of the labels of the expression, hence of $G \setminus \{\varepsilon\}$. In particular, no first includes ε .

Proposition 3. *The expansion of a rational expression is equivalent to the expression.*

Proof. We prove that $\llbracket d(\mathbf{E}) \rrbracket = \llbracket \mathbf{E} \rrbracket$ by induction on the expression. The equivalence is straightforward for (5) to (7) and (10), viz., $\llbracket d(\mathbf{E} \mid \mathbf{F}) \rrbracket = \llbracket d(\mathbf{E}) \mid d(\mathbf{F}) \rrbracket$ (by (10)) = $\llbracket d(\mathbf{E}) \rrbracket \mid \llbracket d(\mathbf{F}) \rrbracket$ (by Lemma 2) = $\llbracket \mathbf{E} \rrbracket \mid \llbracket \mathbf{F} \rrbracket$ (by induction hypothesis) = $\llbracket \mathbf{E} \mid \mathbf{F} \rrbracket$ (by Lemma 2). The case of multiplication, (8), follows from:

$$\begin{aligned} \llbracket d(\mathbf{E} \cdot \mathbf{F}) \rrbracket &= \llbracket d_p(\mathbf{E}) \cdot \mathbf{F} \oplus \langle d_\varepsilon(\mathbf{E}) \rangle \cdot d(\mathbf{F}) \rrbracket &= \llbracket d_p(\mathbf{E}) \rrbracket \cdot \llbracket \mathbf{F} \rrbracket + \langle d_\varepsilon(\mathbf{E}) \rangle \cdot \llbracket d(\mathbf{F}) \rrbracket \\ &= \llbracket d_p(\mathbf{E}) \rrbracket \cdot \llbracket \mathbf{F} \rrbracket + \langle d_\varepsilon(\mathbf{E}) \rangle \cdot \llbracket \mathbf{F} \rrbracket &= \left(\llbracket \langle d_\varepsilon(\mathbf{E}) \rangle \rrbracket + \llbracket d_p(\mathbf{E}) \rrbracket \right) \cdot \llbracket \mathbf{F} \rrbracket \\ &= \llbracket \langle d_\varepsilon(\mathbf{E}) \rangle + d_p(\mathbf{E}) \rrbracket \cdot \llbracket \mathbf{F} \rrbracket &= \llbracket d(\mathbf{E}) \rrbracket \cdot \llbracket \mathbf{F} \rrbracket \\ &= \llbracket \mathbf{E} \rrbracket \cdot \llbracket \mathbf{F} \rrbracket &= \llbracket \mathbf{E} \cdot \mathbf{F} \rrbracket \end{aligned}$$

It might seem more natural to exchange the two terms (i.e., $\langle d_\varepsilon(\mathbf{E}) \rangle \cdot d(\mathbf{F}) \oplus d_p(\mathbf{E}) \cdot \mathbf{F}$), but an implementation first computes $d(\mathbf{E})$ and then computes $d(\mathbf{F})$ only if $d_\varepsilon(\mathbf{E}) \neq 0_{\mathbb{K}}$. The case of Kleene star, (9), follows from Proposition 1. \square

4 Expansion-Based Derived-Term Automaton

Definition 7 (Expansion-Based Derived-Term Automaton). *The derived-term automaton of an expression \mathbf{E} over G is the accessible part of the automaton $\mathcal{A}_{\mathbf{E}} := \langle M, G, \mathbb{K}, Q, E, I, T \rangle$ defined as follows:*

- Q is the set of rational expressions on alphabet A with weights in \mathbb{K} ,
- $I = \mathbf{E} \mapsto 1_{\mathbb{K}}$,
- $E(\mathbf{F}, a, \mathbf{F}') = k$ iff $a \in f(d(\mathbf{F}))$ and $\langle k \rangle \mathbf{F}' \in d(\mathbf{F})(a)$,
- $T(\mathbf{F}) = k$ iff $\langle k \rangle = d(\mathbf{F})(\varepsilon)$.

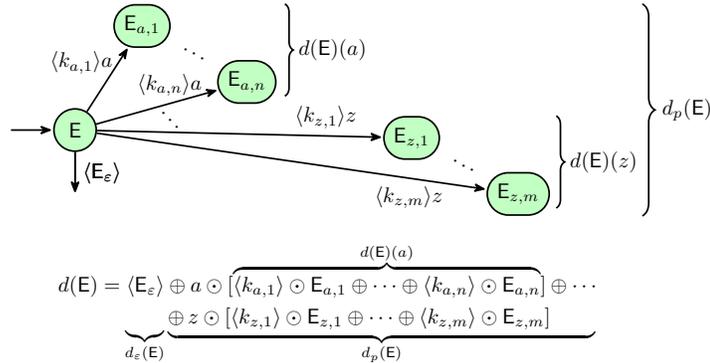


Fig. 2. Initial part of \mathcal{A}_E , the derived-term automaton of E . This figure is somewhat misleading in that some $E_{a,i}$ might be equal to an $E_{z,j}$, or E (but never another $E_{a,j}$).

Since the firsts exclude ε , this automaton is proper. It is straightforward to extract an algorithm from [Definition 7](#), using a work-list of states whose outgoing transitions to compute. The [Fig. 2](#) illustrates the process. This approach admits a natural lazy implementation: the whole automaton is not computed at once, but rather, states and transitions are computed on-the-fly, on demand, for instance when evaluating a word [\[7\]](#). However, we must justify [Definition 7](#) by proving that this automaton is finite ([Theorem 1](#)).

Example 3 ([Examples 1 and 2](#) continued). With $E_1 := \langle 5 \rangle 1 | 1 + \langle 4 \rangle a d e^* | x + \langle 3 \rangle b d e^* | x + \langle 2 \rangle a c e^* | x y + \langle 6 \rangle b c e^* | x y$, one has:

$$\begin{aligned} d(E_1) &= \langle 5 \rangle \oplus a | x \odot [\langle 2 \rangle \odot c e^* | y \oplus \langle 4 \rangle \odot d e^* | \varepsilon] \oplus b | x \odot [\langle 6 \rangle \odot c e^* | y \oplus \langle 3 \rangle \odot d e^* | \varepsilon] \\ &= X_1 \quad (\text{from Example 2}) \end{aligned}$$

[Fig. 1](#) shows the resulting derived-term automaton.

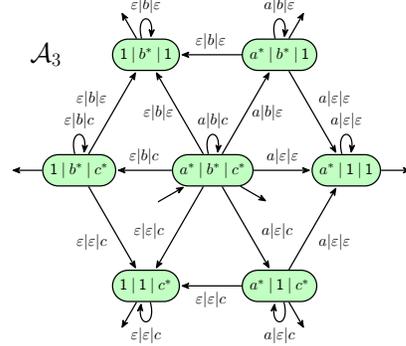
Theorem 1. For any k -tape expression E , $|\mathcal{A}_E| \leq \prod_{i \in [k]} (\|E\|_i + 1) + 1$.

Proof. The proof goes in several steps. First introduce the *true derived terms* of E , a set of expressions noted $\text{TD}(E)$, and the *derived terms* of E , $D(E) := \text{TD}(E) \cup \{E\}$. $\text{TD}(E)$ admits a simple inductive definition similar to [\[2, Def. 3\]](#), to which we add $\text{TD}(E | F) := (\text{TD}(E) | \text{TD}(F)) \cup (\{1\} | \text{TD}(F)) \cup (\text{TD}(E) | \{1\})$, where for two sets of expressions E, F we introduce $E | F := \{E | F\}_{(E,F) \in E \times F}$. Second, verify that $|\text{TD}(E)| \leq \prod_{i \in [k]} (\|E\|_i + 1)$ (hence finite). Third, prove that $D(E)$ is “stable by expansion”, i.e., $\forall F \in D(E), \text{exprs}(d(F)) \subseteq D(E)$. Finally, observe that the states of \mathcal{A}_E are therefore members of $D(E)$, whose size is less than or equal to $1 + |\text{TD}(E)|$. \square

Theorem 2. Any expression E and its expansion-based derived-term automaton \mathcal{A}_E denote the same series, i.e., $\llbracket \mathcal{A}_E \rrbracket = \llbracket E \rrbracket$.

Example 4. Let \mathcal{A}_k be the derived-term automaton of the k -tape expression $a_1^* | \dots | a_k^*$. The states of \mathcal{A}_k are all the possible expressions where the tape i features 1 or a_i^* , except $1 | \dots | 1$. Therefore $|\mathcal{A}_k| = 2^k - 1$, and $\prod_{i \in [k]} (|\mathbb{E}|_i + 1) = 2^k$.

\mathcal{A}_3 , the derived-term automaton of $a^* | b^* | c^*$, is depicted on the right.



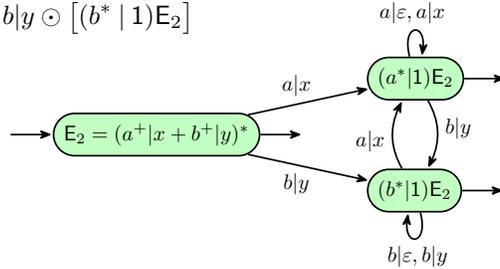
Proof (Theorem 2). We will prove $\llbracket \mathcal{A}_{\mathbb{E}} \rrbracket(m) = \llbracket \mathbb{E} \rrbracket(m)$ by induction on $m \in M$. If $m = \varepsilon$, then $\llbracket \mathcal{A}_{\mathbb{E}} \rrbracket(m) = E_\varepsilon = d(\mathbb{E})(\varepsilon) = \llbracket d(\mathbb{E}) \rrbracket(\varepsilon) = \llbracket \mathbb{E} \rrbracket(\varepsilon)$.

If m is not ε , then it can be generated in a (finite) number of ways: let $F(\mathbb{E}, m) := \{(a, m_a) \in f(d(\mathbb{E})) \times M \mid m = am_a\}$. $F(\mathbb{E}, m)$ is a function: for a given a , there is at most one m_a such that $(a, m_a) \in F(\mathbb{E}, m)$. Fig. 2 is helpful.

$$\begin{aligned}
 \llbracket \mathcal{A}_{\mathbb{E}} \rrbracket(m) &= \sum_{(a, m_a) \in F(\mathbb{E}, m)} \sum_{i \in [n_a]} \langle k_{a,i} \rangle \llbracket \mathcal{A}_{E_{a,i}} \rrbracket(m_a) && \text{by definition of } \mathcal{A}_{\mathbb{E}} \\
 &= \sum_{(a, m_a) \in F(\mathbb{E}, m)} \sum_{i \in [n_a]} \langle k_{a,i} \rangle \llbracket E_{a,i} \rrbracket(m_a) && \text{by induction hypothesis} \\
 &= \sum_{(a, m_a) \in F(\mathbb{E}, m)} \left[\sum_{i \in [n_a]} \langle k_{a,i} \rangle E_{a,i} \right] (m_a) && \text{by Lemma 1} \\
 &= \sum_{(a, m_a) \in F(\mathbb{E}, m)} \llbracket d(\mathbb{E})(a) \rrbracket(m_a) = \sum_{(a, m_a) \in F(\mathbb{E}, m)} \llbracket a \odot d(\mathbb{E})(a) \rrbracket(am_a) \\
 &= \sum_{a \in f(d(\mathbb{E}))} \llbracket a \odot d(\mathbb{E})(a) \rrbracket(m) && F(\mathbb{E}, m) \text{ is a function} \\
 &= \left[\sum_{a \in f(d(\mathbb{E}))} a \odot d(\mathbb{E})(a) \right] (m) && \text{by Lemma 2} \\
 &= \llbracket d_\varepsilon(\mathbb{E}) \rrbracket(m) && \text{by definition} \\
 &= \llbracket d(\mathbb{E}) \rrbracket(m) && \text{since } m \neq \varepsilon \\
 &= \llbracket \mathbb{E} \rrbracket(m) && \text{by Proposition 3} \quad \square
 \end{aligned}$$

Example 5. Let $E_2 := (a^+ | x + b^+ | y)^*$, where $E^+ := EE^*$. Its expansion is $d(E_2) = \langle 1 \rangle \oplus a|x \odot [(a^* | 1)(a^+ | x + b^+ | y)^*] \oplus b|y \odot [(b^* | 1)(a^+ | x + b^+ | y)^*]$
 $= \langle 1 \rangle \oplus a|x \odot [(a^* | 1)E_2] \oplus b|y \odot [(b^* | 1)E_2]$

Its derived-term automaton is:



5 Related Work

Multitape rational expressions have been considered early [11], but “an n-way regular expression is simply a regular expression whose terms are n-tuples of alphabetic symbols or ε ” [9]. However, Kaplan and Kay [9] do consider the full generality of the semantics of operations on rational languages and rational relations, including \times , the Cartesian product of languages, and even use rational expressions more general than their definition. They do not, however, provide an explicit automaton construction algorithm, apparently relying on the simple inductive construction (using the Cartesian product between automata). Our $|$ operator on series was defined as the *tensor product*, denoted \otimes , by Sakarovitch [14, Sec. III.3.2.5], but without equivalent for expressions.

Brzozowski [4] introduced the idea of derivatives of expressions as a means to construct an equivalent automaton. The method applies to extended (unweighted) rational expressions, and constructs a deterministic automaton. Antimirov [3] modified the computation to rely on parts of the derivatives (“partial derivatives”), which results in nondeterministic automata. Lombardy and Sakarovitch [10] extended this approach to support weighted expressions; independently, and with completely different foundations, Rutten [13] proposed a similar construction. Caron et al. [5] introduced support for (unweighted) extended expressions. Demaille [7] provides support for weighted extended expressions; expansions, originally mentioned by Brzozowski [4], are placed at the center of the construct, replacing derivatives, to gain independence with respect to the size of the alphabet, and efficiency. However, the proofs still relied on derivatives, contrary to the present work.

Based on (10) one could attempt to define a derivative-based version, with $\partial_{a|b}(E | F) := \partial_a E | \partial_b F$, however this is troublesome on several regards. First, it would also require $\partial_{a|\varepsilon}$ and $\partial_{\varepsilon|b}$, whose semantics is dubious. Second, from an implementation point of view, that would lead to repeated computations of $\partial_a E$ and of $\partial_b F$, unless one would cache them, but that’s exactly what expansions do. And finally observe that in the derived-term automaton in Example 5, the state $(a^* | 1)(a^+ | x + b^+ | y)^*$ accepts words starting with a on the first tape, and y on the second, yet an outgoing transition on $a|y$ would result in a more complex automaton. Alternative definitions of derivatives may exist², but anyway they would no longer be equivalent to taking the left-quotient of the corresponding language: $a|y$ is a viable prefix from this state.

Different constructions of the derived-term automaton have been discovered [1, 6]. They do not rely on derivatives at all. It is an open question whether these approaches can be adapted to support a tuple operator.

² Makarevskii and Stotskaya [11] define derivatives, but (i) in the case of expressions over tuples of letters, and (ii) only when in so-called “standard form”, for which he notes “no method of constructing [an] n-expression in standard form for a regular n-expression is known.”

6 Conclusion

Our work is in the continuation of derivative-based computations of the derived-term automaton [3–5, 10]. However, we replaced the derivatives by expansions, which lifted the requirement for the monoid of labels to be free.

In order to support k -tape (weighted) rational expressions, we introduced a tupling operator, which is more compact and readable than simple expressions on k -tape letters. We demonstrated how to build the derived-term automaton for any such expressions.

Vcsn¹ implements the techniques exposed in this paper. Our future work aims at other operators, and studying more closely the complexity of the algorithm.

Acknowledgments The author thanks the anonymous reviewers for their constructive comments, and A. Duret-Lutz, S. Lombardy, L. Saiu and J. Sakarovitch for their feedback during this work.

References

1. C. Allauzen and M. Mohri. A unified construction of the Glushkov, follow, and Antimirov automata. In *MFCS*, vol. 4162 of *LNCS*, pp. 110–121. Springer, 2006.
2. P.-Y. Angrand, S. Lombardy, and J. Sakarovitch. On the number of broken derived terms of a rational expression. *Journal of Automata, Languages and Combinatorics*, 15(1/2):27–51, 2010.
3. V. Antimirov. Partial derivatives of regular expressions and finite automaton constructions. *TCS*, 155(2):291–319, 1996.
4. J. A. Brzozowski. Derivatives of regular expressions. *J. ACM*, 11(4):481–494, 1964.
5. P. Caron, J.-M. Champarnaud, and L. Mignot. Partial derivatives of an extended regular expression. In *LATA*, vol. 6638 of *LNCS*, pp. 179–191. Springer, 2011.
6. J.-M. Champarnaud, F. Ouardi, and D. Ziadi. An efficient computation of the equation \mathbb{K} -automaton of a regular \mathbb{K} -expression. In *DLT*, vol. 4588 of *LNCS*. Springer, 2007.
7. A. Demaille. Derived-term automata for extended weighted rational expressions. Technical Report 1605.01530, arXiv, May 2016. URL <http://arxiv.org/abs/1605.01530>.
8. A. Demaille, A. Duret-Lutz, S. Lombardy, and J. Sakarovitch. Implementation concepts in Vaucanson 2. In *CIAA'13*, vol. 7982 of *LNCS*, pp. 122–133, July 2013.
9. R. M. Kaplan and M. Kay. Regular models of phonological rule systems. *Comput. Linguist.*, 20(3):331–378, Sept. 1994.
10. S. Lombardy and J. Sakarovitch. Derivatives of rational expressions with multiplicity. *TCS*, 332(1-3):141–177, 2005.
11. A. Y. Makarevskii and E. D. Stotskaya. Representability in deterministic multi-tape automata. *Cybernetics and System Analysis*, 5(4):390–399, 1969.
12. S. Owens, J. Reppy, and A. Turon. Regular-expression derivatives re-examined. *J. Funct. Program.*, 19(2):173–190, Mar. 2009.
13. J. J. M. M. Rutten. Behavioural differential equations: a coinductive calculus of streams, automata, and power series. *TCS*, 308(1-3):1–53, 2003.
14. J. Sakarovitch. *Elements of Automata Theory*. Cambridge University Press, 2009. Corrected English translation of *Éléments de théorie des automates*, Vuibert, 2003.
15. K. Thompson. Programming techniques: Regular expression search algorithm. *Commun. ACM*, 11(6):419–422, 1968.