



# From regular expressions to deterministic automata

G rard Berry, Ravi Sethi

► **To cite this version:**

G rard Berry, Ravi Sethi. From regular expressions to deterministic automata. [Research Report] RR-0649, INRIA. 1987. <inria-00075904>

**HAL Id: inria-00075904**

**<https://hal.inria.fr/inria-00075904>**

Submitted on 24 May 2006

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destin e au d p t et   la diffusion de documents scientifiques de niveau recherche, publi s ou non,  manant des  tablissements d'enseignement et de recherche fran ais ou  trangers, des laboratoires publics ou priv s.

# INRIA

UNITÉ DE RECHERCHE  
INRIA-SOPHIA ANTIPOLIS

Institut National  
de Recherche  
en Informatique  
et en Automatique

Domaine de Voluceau  
Rocquencourt  
B.P. 105  
78153 Le Chesnay Cedex  
France

Tel. (1) 39 63 55 11

## Rapports de Recherche

N° 649

### FROM REGULAR EXPRESSIONS TO DETERMINISTIC AUTOMATA

Gérard BERRY  
Ravi SETHI

Mars 1987

**FROM REGULAR EXPRESSIONS TO DETERMINISTIC  
AUTOMATA**

**SUR LA TRADUCTION DES EXPRESSIONS REGULIERES EN  
AUTOMATES DETERMINISTES**

*Gérard BERRY \**

*Ravi SETHI\*\**

\* Ecole Nationale Supérieure des Mines de Paris  
Centre de Mathématiques Appliquées  
Valbonne, FRANCE

\*\* BELL Laboratories  
Murray Hill, U.S.A.

# From Regular Expressions to Deterministic Automata

GERARD BERRY

*Ecole des Mines de Paris / INRIA  
Sophia-Antipolis  
06560 Valbonne FRANCE*

RAVI SETHI

*AT&T Bell Laboratories  
Murray Hill, New Jersey 07974 USA*

## Résumé

Notre résultat principal établit une relation entre deux algorithmes de transformation des expressions régulières en automates finis. Le premier, élégant, utilise les "dérivées" d'expressions régulières. Le second, efficace, est fondé sur un "marquage" des expressions.

Les dérivées correspondent directement aux transitions d'états des automates. Quand un automate effectue une transition en lisant une lettre  $a$ , cet  $a$  est enlevé du mot d'entrée. De même, si le mot d'entrée est engendré par une expression  $E$ , alors la dérivée de  $E$  par  $a$  engendre le suffixe de  $a$  du mot d'entrée. L'élégant algorithme de Brzozowski utilise les dérivées pour construire un automate fini: l'état correspondant à une expression  $E$  a une transition étiquetée  $a$  vers l'expression dérivée de  $E$  par  $a$ . Cette approche s'applique à des expressions régulières étendues, par exemple par intersection et négation. Cependant le calcul explicite des dérivées peut être coûteux en pratique.

Le marquage d'expressions, introduit par McNaughton et Yamada, produit des expressions ayant des lettres distinctes. Pour cela on attache des étiquettes aux lettres, transformant par exemple  $(ab + b) * ba$  en  $(a_1b_2 + b_3) * b_4a_5$ . L'algorithme efficace construit d'abord un automate déterministe pour l'expression marquée, puis détermine l'automate non-déterministe obtenu en effaçant les marques. Cette approche ne s'applique cependant pas aux opérateurs comme l'intersection et le complémentaire, pour lesquels le marquage ne préserve pas les langages engendrés.

## Abstract

The main theorem allows an elegant algorithm to be refined into an efficient one. The elegant algorithm for constructing a finite automaton from a regular expression is based on "derivatives of" regular expressions; the efficient algorithm is based on "marking of" regular expressions.

Derivatives of regular expressions correspond to state transitions in finite automata. When a finite automaton makes a transition under input symbol  $a$ , a leading  $a$  is stripped from the remaining input. Correspondingly, if the input string is generated by a regular expression  $E$ , then the derivative of  $E$  by  $a$  generates the remaining input after a leading  $a$  is stripped. Brzozowski used derivatives to construct finite automata; the state for expression  $E$  has a transition under  $a$  to the state for the derivative of  $E$  by  $a$ . This approach extends to regular expressions with new operators, including intersection and complement; however, explicit computation of derivatives can be expensive.

Marking of regular expressions yields an expression with distinct input symbols. Following McNaughton and Yamada, we attach subscripts to each input symbol in an expression;  $(ab+b)*ba$  becomes  $(a_1b_2+b_3)*b_4a_5$ . Conceptually, the efficient algorithm constructs an automaton for the marked expression. The marks on the transitions are then erased, resulting in a nondeterministic automaton for the original unmarked expression. This approach works for the usual operations of union, concatenation, and iteration; however, intersection and complement cannot be handled because marking and unmarking do not preserve the languages generated by regular expressions with these operators.

# From Regular Expressions to Deterministic Automata

GERARD BERRY

*Ecole Nationale Supérieure des Mines de Paris  
Centre de Mathématiques Appliquées  
Sophia-Antipolis, 06560 Valbonne FRANCE*

RAVI SETHI

*AT&T Bell Laboratories  
Murray Hill, New Jersey 07974 USA*

The main theorem allows an elegant algorithm to be refined into an efficient one. The elegant algorithm for constructing a finite automaton from a regular expression is based on "derivatives of" regular expressions; the efficient algorithm is based on "marking of" regular expressions.

Derivatives of regular expressions correspond to state transitions in finite automata. When a finite automaton makes a transition under input symbol  $a$ , a leading  $a$  is stripped from the remaining input. Correspondingly, if the input string is generated by a regular expression  $E$ , then the *derivative of  $E$  by  $a$*  generates the remaining input after a leading  $a$  is stripped. Brzozowski used derivatives to construct finite automata; the state for expression  $E$  has a transition under  $a$  to the state for the derivative of  $E$  by  $a$ . This approach extends to regular expressions with new operators, including intersection and complement; however, explicit computation of derivatives can be expensive.

Marking of regular expressions yields an expression with distinct input symbols. Following McNaughton and Yamada, we attach subscripts to each input symbol in an expression;  $(ab + b)*ba$  becomes  $(a_1b_2 + b_3)*b_4a_5$ . Conceptually, the efficient algorithm constructs an automaton for the marked expression. The marks on the transitions are then erased, resulting in a nondeterministic automaton for the original unmarked expression. This approach works for the usual operations of union, concatenation, and iteration; however, intersection and complement cannot be handled because marking and unmarking do not preserve the languages generated by regular expressions with these operators.

## 1 INTRODUCTION

We study two well-known algorithms for constructing a finite automaton from a regular expression. An elegant algorithm due to Brzozowski [6] will be developed into an efficient algorithm based on McNaughton and Yamada [11]. Brzozowski's algorithm is easily seen to be correct, it accommodates additional operators like intersection and complement, and it has recently served as a starting point for compiling communicating processes in the Esterel programming language [3] into automata [4]. The efficient algorithm is used in fast pattern matchers like *egrep*, distributed as part of the UNIX† operating system. A brief account of *grep* [14] and its cousins, including *egrep*, appears in [1]. A version of *egrep*'s algorithm is described in Section 3.9 of [2].

The syntax of regular expressions over a set  $\Sigma$  of input symbols is ( $a$  is a typical symbol):

$$E ::= 0 \mid 1 \mid a \mid E + E \mid E \cdot E \mid E^*$$

$L(E)$  denotes the language generated by a regular expression  $E$ .  $L(0)$  is the empty set and  $L(1)$  is the set consisting of the empty string  $\epsilon$ . Note that 0 and 1 are not input symbols; they represent the sets of strings  $\emptyset$  and  $\{\epsilon\}$ .  $L(E + F)$  is the union of  $L(E)$  and  $L(F)$ .  $L(E \cdot F)$  consists of strings formed by concatenating a string in  $L(E)$  with a string in  $L(F)$ .  $L(E^*)$  consists of strings formed by concatenating zero or more strings from  $L(E)$ ;  $L(E^*)$  includes the empty string  $\epsilon$ .

We write  $E \equiv F$  if  $L(E) = L(F)$ . The following properties of regular expressions will be used

† UNIX is a registered trademark of AT&T.

without fanfare:

$$\begin{aligned} 0 + E &\equiv E \equiv E + 0 \\ 0 \cdot E &\equiv 0 \equiv E \cdot 0 \\ 1 \cdot E &\equiv E \equiv E \cdot 1 \end{aligned}$$

Using Brzozowski's notation,  $\delta(E)$  stands for 1 if  $L(E)$  contains the empty string; otherwise,  $\delta(E)$  stands for 0. It can be computed from the structure of  $E$ :

$$\begin{aligned} \delta(0) &= 0 \\ \delta(1) &= 1 \\ \delta(a) &= 0 \\ \delta(E+F) &= \delta(E) + \delta(F) \\ \delta(E \cdot F) &= \delta(E) \cdot \delta(F) \\ \delta(E^*) &= 1 \end{aligned}$$

Thus,  $\delta(E) \cdot F$  equals  $F$  if the empty string is in  $L(E)$ ; otherwise,  $\delta(E) \cdot F$  equals 0. Furthermore,  $E + \delta(E) \equiv E$ , because  $E + 0 \equiv E$  and, when  $\varepsilon$  is in  $L(E)$ ,  $E + 1 \equiv E$ .

## 2 DERIVATIVES OF REGULAR EXPRESSIONS

Brzozowski's algorithm [6] relies on the notion of a "derivative" of a regular expression  $E$  with respect to an input symbol  $a$ , written  $a^{-1}E$ . Informally, if leading  $a$ 's are stripped from strings in  $L(E)$  that start with  $a$ , we get the strings generated by  $a^{-1}E$ . The derivative of  $aba + bb$  by  $a$  is  $ba$ .

DEFINITION 1. Given a regular expression  $E$  and a symbol  $a$ , the *derivative of  $E$  by  $a$* , written  $a^{-1}E$ , is defined by:

$$\begin{aligned} a^{-1}1 &= 0 \\ a^{-1}0 &= 0 \\ a^{-1}a &= 1 \\ a^{-1}b &= 0 && \text{for } b \neq a \\ a^{-1}(E+F) &= a^{-1}E + a^{-1}F \\ a^{-1}(E \cdot F) &= a^{-1}E \cdot F + \delta(E) \cdot a^{-1}F \\ a^{-1}(E^*) &= a^{-1}E \cdot E^* \end{aligned}$$

Within expressions,  $a^{-1}$  is treated as a prefix operator with higher precedence than  $+$ ,  $\cdot$ , and  $*$ .  $\square$

REMARK. Additional operators like intersection and complement can be handled by adding rules of the form:

$$\begin{aligned} a^{-1}(E \cap F) &= a^{-1}E \cap a^{-1}F \\ a^{-1}(E - F) &= a^{-1}E - a^{-1}F \end{aligned}$$

Informally, it does not matter if leading  $a$ 's are stripped before or after the operations are performed. For example,

$$\begin{aligned} a^{-1}(ab^* \cap a) &\equiv 1 \equiv (a^{-1}(ab^*)) \cap (a^{-1}a) \\ a^{-1}(ab^* - a) &\equiv b^* - 1 \equiv (a^{-1}(ab^*)) - (a^{-1}a) \end{aligned} \quad \square$$

States in the constructed automaton correspond to regular expressions. There is a transition under  $a$  from the state for  $E$  to the state for  $a^{-1}E$ . The transition under  $a$  from the state for  $aba + bb$  is to the state for  $ba$ ; the subsequent transition under  $b$  is to the state for  $a$ . It is easier to talk about sequences of transitions if the notion of derivatives is generalized from symbols to strings. It is easier to write  $(ab)^{-1}E$  than  $b^{-1}(a^{-1}E)$ . More significantly, the next definition allows us to write  $w^{-1}E$ , where the string  $w$  is any member of a set of strings.

DEFINITION 2. The extension from symbols to the derivative of  $E$  by a string is defined by:

$$\begin{aligned} \epsilon^{-1}E &= E \\ (wa)^{-1}E &= a^{-1}(w^{-1}E) \end{aligned}$$

Within expressions,  $w^{-1}$  is treated as a prefix operator with higher precedence than  $+$ ,  $\cdot$ , and  $*$ .  $\square$

Automata will be constructed by explicitly computing derivatives, then derivatives of derivatives, and so on, as needed. Convergence is guaranteed by the following result from [6].

PROPOSITION 3. *The set of derivatives of a regular expression is finite, modulo associativity, commutativity, and idempotence of  $+$ ; that is, the set  $\{F \mid w.F = w^{-1}E\}$  has a finite number of equivalence classes.*  $\square$

Without associativity, commutativity, and idempotence of  $+$ , duplicate subexpressions would cause successive derivatives of  $E = a^*(aa)^*$  by  $a, aa, \dots$  to be distinct. These properties allow a sum of expressions to be treated as a set of expressions, thereby removing duplicates.

ALGORITHM 1. ([6]) Construction of a deterministic automaton  $D$  accepting  $L(E)$ .

1. The states of  $D$  are the distinct derivatives  $w^{-1}E$ , for all strings  $w$ . Proposition 3 ensures convergence of this step.
2. Construct a transition under  $a$  from state  $p$  to state  $q$  if and only if  $p$  is for derivative  $w^{-1}E$ , for some  $w$ , and  $q$  is for  $(wa)^{-1}E$ .
3. The state for  $E$  is the start state. A state is an accepting state if and only if it is for a derivative  $w^{-1}E$ , for some  $w$ , and  $\delta(w^{-1}E) = 1$ ; that is, the empty string is in  $L(w^{-1}E)$ .  $\square$

Algorithm 1 constructs the automaton in Fig. 1. Each state of the automaton is for a derivative of  $E = (ab + b)^*ba$ ; the state for  $w^{-1}E$  has a transition under  $a$  to the state for  $(wa)^{-1}E$ .

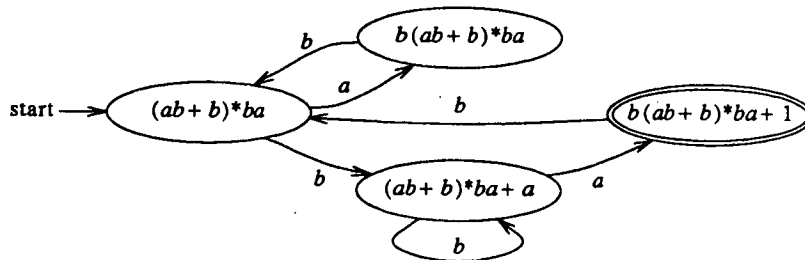


Fig. 1. Automaton accepting  $(ab + b)^*ba$ .

### 3 REGULAR EXPRESSIONS WITH DISTINCT SYMBOLS

Following McNaughton and Yamada [11], we mark all input symbols in a regular expression to make them distinct. The marks are written as subscripts; a marked version of  $(ab + b)^*ba$  is  $(a_1b_2 + b_3)^*b_4a_5$ , where  $a_1$  and  $a_5$  are treated as different symbols.

The construction of a deterministic automaton  $D$  from a regular expression  $E$  is outlined in Fig. 2.  $E'$  is formed by marking all symbols in  $E$  to make them distinct. Suppose that an automaton  $M'$  accepts  $L(E')$ . We show that unmarking the symbols labeling the transitions of  $M'$  yields a nondeterministic automaton  $M$ .  $D$  can be obtained from  $M$  by applying the standard subset construction [13, 2].

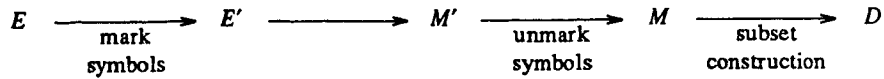


Fig. 2. From a regular expression  $E$  to a deterministic finite automaton  $D$ .

REMARK 4. The approach of Fig. 2 does not extend to regular expressions with intersection and complement operators. Although  $(ab^*) \cap a \equiv a$ , we get  $(a_1b_2^*) \cap a_3 \equiv 0$  because  $a_1$  and  $a_3$  are distinct. Similarly,  $(ab^*) - a \equiv abb^*$ , but  $(a_1b_2^*) - a_3 \equiv a_1b_2^*$ . The unmarking homomorphism does not commute with intersection and relative complementation of languages.

The approach of Fig. 2 works for union, concatenation, and iteration because unmarking commutes with these operations on languages.  $\square$

PROPOSITION 5. Let  $E'$  be the regular expression obtained from  $E$  by marking all symbols to make them distinct. If  $M'$  is an automaton accepting  $L(E')$ , then  $M$ , a nondeterministic automaton accepting  $L(E)$ , is obtained by unmarking all the symbols labeling the transitions in  $M'$ .

*Proof.* Let  $unmark$  map marked symbols to their original unmarked form. We write  $L(A)$  for the language accepted by automaton  $A$ .

A structural induction on  $E$  establishes that  $L(E) = unmark(L(E'))$ . For the basis, note that marking has no effect on 0 and 1, and that unmarking recovers a symbol  $a$  from its marked counterpart. The inductive step, consisting of cases for the operators  $+$ ,  $\cdot$ , and  $*$ , is not shown.

For each transition of  $M'$  on a marked symbol, there is a corresponding nondeterministic transition of  $M$  on the unmarked symbol, so  $M$  accepts a string  $w$  if  $M'$  accepts a string  $w'$  such that  $w = unmark(w')$ . The converse holds as well, so  $L(M) = unmark(L(M'))$ .

By construction,  $L(E') = L(M')$ . The result  $L(E) = L(M)$  follows by transitivity.  $\square$

The main theorem in this section allows each symbol in a marked expression to be viewed as a state of an automaton. Figure 3 contains a motivating example for the theorem. Each state of the automaton in Fig. 3 is labeled with a symbol representing a derivative of the expression  $(a_1b_2 + b_3)^*b_4a_5$ . Furthermore, as in automata constructed by Algorithm 1, the state for  $C$  has a transition under a marked symbol to a state for the derivative of  $C$  by that symbol. By construction, all the transitions entering a state are labeled with the same symbol; see for example the transitions labeled  $b_3$  into the state for  $C_3$ . All strings that drive the automaton from the start state into  $C_3$  must therefore be of the form  $wb_3$ , for some  $w$ . By construction, if a string  $wb_3$  drives the automaton into a state for expression  $C_3$ , then  $C_3$  must equal the derivative by  $wb_3$  of the starting expression  $C_0$ . Theorem 7 shows that the equivalence between  $C_3$  and  $(wb_3)^{-1}C_0$ , for all  $w$ , is no accident; it follows from the distinctness of all symbols in the expression  $C_0$ .

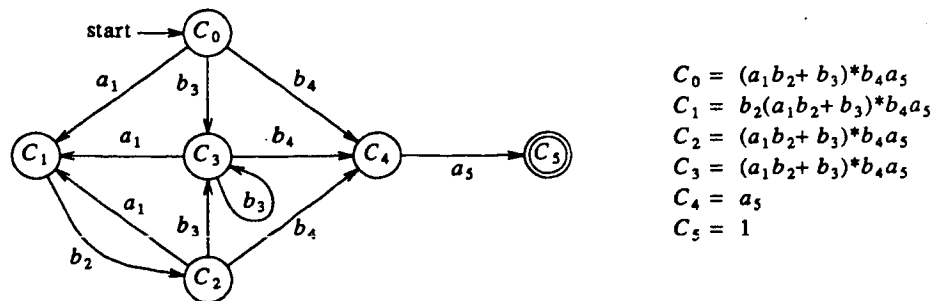


Fig. 3. Automaton for  $(a_1b_2 + b_3)^*b_4a_5$ .

The proof of the main theorem uses the following lemma.



LEMMA 6. Given any symbol  $a$ , for all strings  $w$ ,  $(wa)^{-1}(E^*)$  is equivalent to a sum of subterms chosen from the set  $\{(va)^{-1}E \cdot E^* \mid wa = uva\}$ .

*Proof.* By induction on the length of  $w$ . The basis, length zero, follows from the definition of derivatives because

$$a^{-1}(E^*) = a^{-1}E \cdot E^*$$

For the inductive step, suppose  $w = xb$ . By definition,

$$(xba)^{-1}(E^*) = a^{-1}((xb)^{-1}(E^*))$$

From the inductive hypothesis,  $(xb)^{-1}(E^*)$  is equivalent to a sum of subterms chosen from  $\{(zb)^{-1}E \cdot E^* \mid xb = yzb\}$ . The operator  $a^{-1}$  distributes across a sum:

$$a^{-1}\left(\sum_{x=yz} (zb)^{-1}E \cdot E^*\right) = \sum_{x=yz} a^{-1}((zb)^{-1}E \cdot E^*)$$

Applying the rule for  $\cdot$  in Definition 1, the right side yields

$$\sum_{x=yz} (zba)^{-1}E \cdot E^* + \sum_{x=yz} \delta((zb)^{-1}E) \cdot a^{-1}(E^*) \quad (1)$$

Each subterm of the form  $\delta((zb)^{-1}E)$  is either 0 or 1. Therefore, the second summation in (1) equals  $a^{-1}(E^*) = a^{-1}E \cdot E^*$ , if it does not equal 0. All terms in (1) are therefore of the form  $(va)^{-1}E \cdot E^*$ , where  $wa = uva = xba$ , so the lemma holds.  $\square$

THEOREM 7 (Main Theorem). Let all symbols in  $E$  be distinct. Given any symbol  $a$ , for all strings  $w$ ,  $(wa)^{-1}E$  is either 0 or unique modulo associativity, commutativity, and idempotence of  $+$ .

*Proof.* By structural induction on  $E$ . If  $E$  is either 0 or 1, then all derivatives are 0. Otherwise, if it is a symbol, then the only possible values for its derivatives are 0 and 1. In the remaining cases, we use the structure of  $E$  to expand  $(wa)^{-1}E$ .

*Case:*  $E = E_1 + E_2$ . It follows from the definition of derivatives that

$$(wa)^{-1}(E_1 + E_2) = (wa)^{-1}E_1 + (wa)^{-1}E_2 \quad (2)$$

Since all symbols in  $E$  are distinct, if  $a$  is in  $E_1$  then  $(wa)^{-1}E_2 = 0$ ; otherwise,  $a$  is in  $E_2$  and  $(wa)^{-1}E_1 = 0$ . This case follows from the inductive hypothesis applied to the remaining term in (2).

*Case:*  $E = E_1 \cdot E_2$ . An auxiliary induction on the length of  $w$  establishes

$$(wa)^{-1}(E_1 \cdot E_2) = (wa)^{-1}E_1 \cdot E_2 + \sum_{wa=uv a} \delta(u^{-1}E_1) \cdot (va)^{-1}E_2 \quad (3)$$

If  $a$  is in  $E_1$ , then only the first term on the right side of equality (3) survives and the inductive hypothesis applies to it. Otherwise,  $a$  is in  $E_2$  and terms of the form  $\delta(u^{-1}E_1) \cdot (va)^{-1}E_2$  are left. Recall that  $\delta(u^{-1}E_1)$  is either 0 or 1, for all  $u$ , and from the inductive hypothesis, all subterms of the form  $(va)^{-1}E_2$  are equivalent to either 0 or some fixed term  $F$ , so their sum is also either 0 or equivalent to that fixed term  $F$ .

*Case:*  $E = E_1^*$ . From Lemma 6,  $(wa)^{-1}(E^*)$  is equivalent to a sum of terms chosen from the set  $\{(va)^{-1}E \cdot E^* \mid wa = uva\}$ ; a complete characterization of  $(wa)^{-1}(E^*)$  is not needed. From the inductive hypothesis, each nonzero subterm  $(va)^{-1}E$  must therefore be equivalent to some fixed term  $F$ , so any sum of subterms of the form  $(va)^{-1}E \cdot E^*$  is equivalent to  $F \cdot E^*$ . This final case holds because, being such a sum,  $(wa)^{-1}(E^*)$  is equivalent to  $F \cdot E^*$  if it is not 0.  $\square$

The automaton in Fig. 3 is constructed by a refinement of Algorithm 1. The new algorithm works with a specific set of derivatives, corresponding to the distinct symbols in a marked expression. Since this set is fixed in advance, convergence of the new algorithm is immediate; correctness carries over from Algorithm 1.

DEFINITION 8. Let all symbols in  $E$  be distinct. For all symbols  $a$  in  $E$ , a *continuation* of  $a$  in  $E$  is

any expression  $(wa)^{-1}E \neq 0$ . By structural induction on  $E$ , such an expression must exist; by the above theorem, all such expressions are equivalent. We therefore speak of "the" continuation of  $a$  in  $E$  to refer to some expression in the equivalence class.  $\square$

ALGORITHM 2. Construction of a deterministic automaton  $M'$  from a marked expression  $E'$ .

1.  $M'$  has a state for the continuation of each marked symbol in  $E'$ .
2. Construct a transition under  $a$  from state  $p$  to the state for the continuation of  $a$  if and only if  $p$  is for some continuation  $C$  and  $C$  can generate a string with a leading  $a$ .
3. The state for the entire expression  $E'$  is the start state. A state is an accepting state if and only if it is for a continuation  $C$  and  $\delta(C) = 1$ .  $\square$

The states of the automaton in Fig. 3 are labeled with continuations of marked input symbols in the expression  $(a_1b_2 + b_3)^*b_4a_5$ . For all  $i$ ,  $1 \leq i \leq 5$ ,  $C_i$  is the continuation of the symbol marked  $i$ , and  $C_0$  represents the entire expression. Although the automata constructed by the new algorithm can have more states — compare Fig. 3 and Fig. 1 — the new algorithm compensates by not checking expressions for equivalence. The next section shows that the continuations in Fig. 3 need not be computed explicitly either.

#### 4 A FAST ALGORITHM

Algorithm 2 can be improved. Since each continuation is uniquely determined by an input symbol, a regular expression with  $n$  marked symbols will lead to an automaton with  $n+1$  states — a start state and a state for each symbol. Instead of complete continuations, the second step of the algorithm needs only the set of leading symbols in strings generated by the continuations. These sets are related to "follow sets" defined below.

DEFINITION 9.

$$\begin{aligned} \text{first}(E) &= \{a \mid av \in L(E)\} \\ \text{follow}_E(a) &= \{b \mid uabv \in L(E)\} \end{aligned} \quad \square$$

Expressions of the form  $E!$ , where  $!$  is a new endmarker symbol, are used below to avoid special cases in the computation of follow sets for the "last" symbols that can be generated by  $E$ . If  $a$  is such a last symbol, then  $\text{follow}_{E!}(a)$  will contain  $!$ .

PROPOSITION 10. *Let all symbols in  $E$  be distinct and for all  $a$  let  $C_a$  be the continuation of  $a$  in  $E$ . Then,*

$$a. \text{first}(C_a!) = \text{follow}_{E!}(a)$$

*Proof.* Brzozowski's [6] observation that every regular expression  $E$  can be represented as an infinite sum of terms of the form  $w \cdot w^{-1}E$ , formalizes the idea that the derivative of  $w^{-1}E$  is formed by stripping a prefix  $w$  from strings generated by  $E$ . Restating the result for the derivatives of  $E$  by strings of the form  $wa$

$$E \equiv \delta(E) + \sum_{wa} wa \cdot (wa)^{-1}E$$

From Theorem 7,  $(wa)^{-1}E$  is either 0 or  $C_a$ , so

$$E! \equiv \delta(E) \cdot ! + \sum_{(wa)^{-1}E \neq 0} wa \cdot C_a \cdot ! \quad (4)$$

For all  $a$ , symbol  $a$  appears just before  $C_a$  in the right side of (4), so  $\text{first}(C_a!)$  is a subset of  $\text{follow}_{E!}(a)$ .

For the converse, suppose  $b$  is in  $\text{follow}_{E!}(a)$ . Then, for some  $u$  and  $v$ ,  $uabv! \in L(E!)$ , and from (4)

$$uabv! \in L(ua \cdot C_a \cdot !)$$

Hence,  $b$  must be in  $first(C_a!)$ . □

An algorithm for computing follow sets is given in Section 3.9 of [2]. A related algorithm is illustrated in Fig. 4. The figure contains a syntax tree for the expression  $(a_1b_2 + b_3)^*b_4a_5!$ , where  $!$  is an endmarker. For clarity, follow sets are written in terms of integer subscripts on the marked symbols. The notion of follow sets is generalized from symbols to subexpressions represented by nodes in the syntax tree. To the right of each node appears the set of symbols that can follow the last symbol generated by the subexpression at the node. Alternatively, the set can be thought of as an attribute inherited by the node. Proceeding top down from the root with  $\emptyset$ , the set at a node is accumulated until the leaf for a symbol  $a$  is reached with  $follow_{E_1}(a)$ .

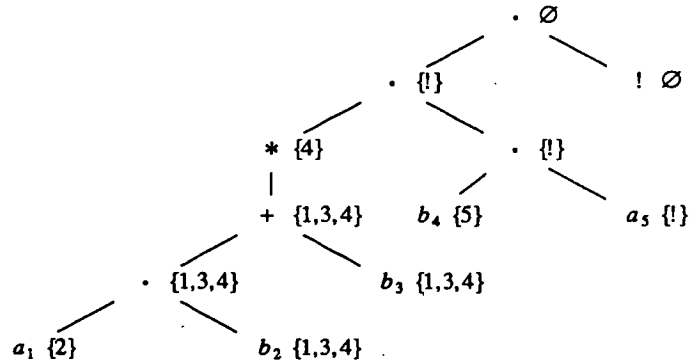


Fig. 4. Follow sets for subexpressions of  $(a_1b_2 + b_3)^*b_4a_5!$ .

The rules used for computing the sets in Fig. 4 are summarized in the next proposition.

PROPOSITION 11. Let  $E$  be a regular expression with distinct symbols.  $F$ , defined by the rules below, is such that  $F(E, \{!\})$  yields a set of pairs of the form  $\langle a, follow_{E_1}(a) \rangle$ . The rules are:

$$\begin{aligned} F(E_1 + E_2, S) &= F(E_1, S) \cup F(E_2, S) \\ F(E_1 \cdot E_2, S) &= F(E_1, first(E_2 \cup \delta(E_2) \cdot S)) \cup F(E_2, S) \\ F(E_1^*, S) &= F(E_1, first(E_1) \cup S) \\ F(a, S) &= \langle a, S \rangle \\ F(1, S) &= \emptyset \\ F(0, S) &= \emptyset \end{aligned}$$

□

Clément and Kahn [8] have adapted the above rules to construct a Typol [7] implementation. The Typol specification is built up of sequents of the form

$$S \leftarrow E_1 : \langle first, \delta, F \rangle$$

where  $first = first(E_1)$ ,  $\delta = \delta(E_1)$ , and  $F = F(E_1, S)$ . Their rule for  $E_1 \cdot E_2$  is

$$\frac{first(E_2) \cup \delta(E_2) \cdot S \leftarrow E_1 : \langle first_1, \delta_1, F_1 \rangle \quad S \leftarrow E_2 : \langle first_2, \delta_2, F_2 \rangle}{S \leftarrow E_1 \cdot E_2 : \langle first_1 \cup \delta_1 \cdot first_2, \delta_1 \cdot \delta_2, F_1 \cup F_2 \rangle}$$

Once  $first$  and  $follow$  sets are computed, algorithm 2 is improved into algorithm 3 below.

ALGORITHM 3. (Fast construction of a deterministic automaton  $M'$  from a marked expression  $E'$ ).

1.  $M'$  has a start state plus a state for each marked symbol  $a_i$  in  $E'$ .
2. Construct a transition under  $a_i$  from the start state to the state for  $a_i$  if and only if  $a_i \in first(E')$ ; construct a transition under  $a_i$  from the state for  $b_j$  to the state for  $a_i$  if and only if  $a_i \in follow_{E'}(b_j)$ .
3. The start state is an accepting state if and only if  $\delta(E') = 1$ ; the state for  $a_i$  is an accepting state if and only if  $! \in follow_{E'}(a_i)$ .

## 5 DISCUSSION

We no longer need to choose between the simpler and more general approach based on derivatives of regular expressions and the efficient approach based on marking and *follow* sets. Theorem 7 and Proposition 10 relate the two approaches, so the derivative approach can be used as the starting point for constructing automata. These results also verify the correctness of algorithms based on marking and *follow* sets.

Regular expressions with union, concatenation, and iteration operators suffice for specifying patterns in strings [14, 1]. Additional operators are used in applications like protocol validation [5, 9] and communication between processes [12, 3, 10]. The idea of marking and unmarking symbols in regular expressions does not always extend to additional operators; see the examples involving intersection and complement in Remark 4. The idea of *follow* sets can, however, be used even if marking cannot.

## ACKNOWLEDGMENTS

We thank Al Aho, Dominique Clément, and Gilles Kahn for helpful comments.

## REFERENCES

1. AHO, A. V. "Pattern matching in strings," in R. V. Book, *Formal Language Theory*, Academic Press, New York (1980) pp. 325-347.
2. AHO, A. V., R. SETHI, AND J. D. ULLMAN "Compilers: Principles, Techniques, and Tools," Addison-Wesley, Reading, Mass. (1986).
3. BERRY, G. AND L. COSSERAT "The Esterel synchronous programming language and its mathematical semantics," in S. D. Brookes, A. W. Roscoe, and G. Winskel, *Seminar in Concurrency*, Lecture Notes in Computer Science 197, Springer-Verlag, Berlin (1985).
4. BERRY, G., P. COURONNE, AND G. GONTHIER "The ESTEREL v2.1 System Manuals," Tech. Report, Ecole des Minés/INRIA (1986).
5. BOCHMANN, G. V. "Communication protocols and error recovery procedures," *ACM Operating Systems Review* 9:3 (March 1975) 45-50.
6. BRZOZOWSKI, J. A. "Derivatives of regular expressions," *J. ACM* 11:4 (1964) 481-494.
7. CLEMENT, D., J. DESPEYROUX, T. DESPEYROUX, L. HASCOET, AND G. KAHN "Natural semantics on the computer," No. 416, INRIA, Sophia-Antipolis (1985).
8. CLEMENT, D. AND G. KAHN "Personal communication (February 1986).
9. HOLZMANN, G. J. "A theory for protocol validation," *IEEE Trans. Computers* C-31:8 (1982) 730-738.
10. KATZENELSON, J. AND R. P. KURSHAN "S/R: A language for specifying protocols and other coordinating processes," *5th Phoenix Conference on Computer Communications* (1986) 286-292.
11. McNAUGHTON, R. AND H. YAMADA "Regular expressions and state graphs for automata," *IRE Trans. on Electronic Computers* EC-9:1 (1960) 38-47.
12. MILNER, R. "A complete inference system for a class of regular behaviours," *J. Computer and System Sciences* 28 (1984) 439-466.
13. RABIN, M. O. AND D. SCOTT. "Finite automata and their decision problems," *IBM J. Research and Development* 3:2 (1959) 114-125.
14. THOMPSON, K. "Regular expression search algorithm," *Comm. ACM* 11:6 (1968) 419-422.

Imprimé en France

par

l'Institut National de Recherche en Informatique et en Automatique

