

# Regular expressions into finite automata

Anne Brüggemann-Klein

*Institut für Informatik, Universität Freiburg, Rheinstr. 10–12, 7800 Freiburg, Germany*

Communicated by M. Nivat  
Received September 1991  
Revised June 1992

## Abstract

Brüggemann-Klein, A., Regular expressions into finite automata, Theoretical Computer Science 120 (1993) 197–213.

It is a well-established fact that each regular expression can be transformed into a nondeterministic finite automaton (NFA) with or without  $\varepsilon$ -transitions, and all authors seem to provide their own variant of the construction. Of these, Berry and Sethi (1986) have shown that the construction of an  $\varepsilon$ -free NFA due to Glushkov (1961) is a *natural* representation of the regular expression because it can be described in terms of the Brzozowski derivatives (Brzozowski 1964) of the expression. Moreover, the Glushkov construction also plays a significant role in the document processing area: The SGML standard (ISO 8879 1986), now widely adopted by publishing houses and government agencies for the syntactic specification of textual markup systems, uses *deterministic* regular expressions, i.e. expressions whose Glushkov automaton is deterministic, as a description language for document types.

In this paper, we first show that the Glushkov automaton can be constructed in a time quadratic in the size of the expression, and that this is worst-case optimal. For deterministic expressions, our algorithm has even linear run time. This improves on the cubic time methods suggested in the literature (Book et al. 1971; Aho et al. 1986; Berry and Sethi 1986). A major step of the algorithm consists in bringing the expression into what we call *star normal form*. This concept is also useful for characterizing the relationship between two types of unambiguity that have been studied in the literature. Namely, we show that, modulo a technical condition, an expression is strongly unambiguous (Sippu and Soisalon-Soininen 1988) if and only if it is weakly unambiguous (Book et al. 1971) and in star-normal form. This leads to our third result, a quadratic-time decision algorithm for weak unambiguity, that improves on the biquadratic method introduced by Book et al. (1971).

## 1. Introduction

Regular expressions play a prominent role in practical applications. In syntactic specifications of programming languages they describe lexical tokens, and in text

*Correspondence to:* A. Brüggemann-Klein, Institut für Informatik, Universität Freiburg, Rheinstr. 10–12, 7800 Freiburg, Germany.

manipulation systems they describe textual patterns that trigger processing actions [2, 9]. They have become the basis of standard utilities such as the scanner generator *lex* and the text tools *awk* and *egrep* [1, 16]. Regular expressions provide an appropriate notation for regular languages in text-based user interfaces, whereas finite automata are the preferred internal data structure for programming purposes.

Two distinct methods have been devised to translate a regular expression into a nondeterministic finite automaton (NFA). In a two-step approach, the standard method first translates a regular expression in linear time into a nondeterministic finite automaton with  $\varepsilon$ -transitions. Then the  $\varepsilon$ -transitions are eliminated in quadratic time [12, 2, 17, 15].

The alternative method formalizes the notion of a symbol in a word being matched by an occurrence of the symbol in the expression [10, 5, 2, 4]. It is based on the fact that, if a word is denoted by an expression, it must be possible to spell out that word by tracing an appropriate “path” through the expression. For example, the word *abba* is denoted by the expression  $a(a+b)^*a$  because it corresponds to the path that starts at the first *a* in  $a(a+b)^*a$ , then visits the *b* twice and finally arrives at the third occurrence of *a*. Of course, the structure of the expression restricts the positions that adjacent symbols of a word can be matched with. For instance, if the  $i$ th symbol in a word is matched by the second *a* in  $a(a+b)^*a$ , then the  $(i+1)$ th symbol cannot be matched with the first *a*. These restrictions were first formalized by Glushkov [10].

It has been noted by a number of authors that a regular expression  $E$  defines in a natural way an NFA  $M_E$ , the *Glushkov automaton* of  $E$ , whose states correspond to the occurrences of symbols in  $E$  and whose transitions connect positions that can be consecutive on a path through  $E$  [5, 2]. Recently, Berry and Sethi have shown that the Glushkov construction  $M_E$  is related in a natural way to the Brzowski derivatives of  $E$  [7, 4].

None of the cited papers considers, however, the time complexity of constructing  $M_E$ . A straightforward implementation takes time cubic in the size of the expression, as opposed to the quadratic time of the standard construction.

In this paper we provide a quadratic-time algorithm (Theorem 3.9) that is worst-case optimal and output sensitive. To this end, we first transform an expression  $E$  in linear time into an expression  $E^*$  in what we denote by *star-normal form* whose Glushkov automaton is identical to  $M_E$  (Theorem 3.1). Then we show how, for expressions in star normal form, the Glushkov automaton can be constructed in quadratic time (Lemma 2.7). An alternative proof for Theorem 3.9 has independently been found by Chen and Paige [8].

For some practical applications, full regular expressions are considered too powerful and syntactic restrictions are imposed. Well-known examples are the Unix tools *vi* and *grep* [9]. In the text processing area, the ISO Standard for SGML (standard general markup language) provides a syntactic metalanguage for the definition of textual markup systems. Such markup systems facilitate the electronic interchange of electronic documents and provide a standard basis for accessing and displaying them.

In the SGML context, the only valid regular expressions are those for which the Glushkov automaton is deterministic. The languages recognized by *deterministic* regular expressions have been characterized [6]. Here we show that for a deterministic expression a *deterministic* finite automaton can be constructed in linear time. This implies that LL(1) parsing tables of linear size can be generated for the context-free grammars SGML uses to describe document types.

When transforming language descriptions from one type to another, e.g. from regular expressions to finite automata, it is, from an applications point of view, important to preserve unambiguity, since only for unambiguous representations of a language can the meaning of a word in the language be derived from the representation. Indeed, this was the motivation for Book et al. [5] to investigate the NFA  $M_E$ . They showed that a regular expression  $E$  is unambiguous if and only if  $M_E$  is unambiguous.

An  $\varepsilon$ -NFA  $M$  is *unambiguous* if for each work  $w$ , there is at most one path through the state diagram of  $M$  that spells out  $w$  [2]. A regular expression  $E$  is *unambiguous* if, for each word  $w$ , there is at most one path through  $E$  that matches  $w$  [5]. Thus, in unambiguous  $\varepsilon$ -NFAs, semantic procedures can be attached to transitions, and in unambiguous regular expressions, they can be attached to occurrences of symbols.

We call the kind of ambiguity for regular expressions as defined above *weak*, as opposed to another definition given by Sippu and Soisalon-Soininen [15]. Their *strong* unambiguity allows semantic procedures to be attached not only to the symbols but also to the operators in a regular expression. To give an example, the expression  $(a^* + b^*)^*$  is trivially weakly unambiguous because each symbol occurs only once. Thus, any symbol in a word can be matched by exactly one position in the expression. In contrast, the word  $aa$  is denoted by  $(a^* + b^*)^*$  as a single application of the outer star and a twofold application of the inner one or, alternatively, as a twofold application of the outer star and two single applications of the inner one. Thus,  $(a^* + b^*)^*$  is *not* strongly unambiguous.

The two notions of unambiguity are related via our notion of star normal form. In Theorem 4.9 we show that, essentially, an expression is strongly unambiguous if and only if it is weakly unambiguous and in star normal form.

Finally, we turn to the decision problem for weak unambiguity. Unambiguity of  $\varepsilon$ -NFAs can be reduced in linear time to the LR(0) property for context-free grammars, which has quadratic-time complexity [15]. Strong unambiguity of expressions can be reduced in linear time to unambiguity of  $\varepsilon$ -NFAs via Thompson's construction [2]. Thus, there is a quadratic-time algorithm to decide whether an expression is strongly unambiguous. On the other hand, weak unambiguity can as well be reduced to unambiguity of NFAs via the Glushkov construction, but because the reduction is quadratic in time and size, this yields a biquadratic decision algorithm for weak unambiguity. Alternatively, Book et al. suggest testing a regular expression  $E$  for weak ambiguity by transforming  $M_E$  into a Mealy automaton that is then tested for information losslessness. The latter can be done using an algorithm by Huffman or

Evans, given, for example, in the textbook of Hennie [11]. This algorithm boils down to testing, for any two different states  $p$  and  $q$  of  $M_E$  that can be reached from the initial state by means of a common word  $w$ , whether there is a state  $r$  and transitions from  $p$  to  $r$  and  $q$  to  $r$  on a common symbol  $a$ . Essentially,  $M_E$  and, hence,  $E$  is (weakly) unambiguous if no such pair of states can be found. A straightforward implementation of this algorithm is biquadratic in the size of  $E$ , too.

Applying the technique developed for the quadratic-time construction of the Glushkov automaton, we can transform a regular expression  $E$  into  $E^*$  in star normal form in linear time. This transformation preserves weak unambiguity and, for expressions in star-normal form, weak and strong unambiguity are essentially the same. Thus, we provide the first algorithm for deciding weak unambiguity in quadratic time (Theorem 4.11).

## 2. Definitions

In this section, we define the Glushkov NFA  $M_E$  for a regular expression  $E$ . A straightforward implementation of the construction runs in time cubic in the size of  $E$ . We show that the implementation can be modified to run in quadratic time, provided that  $E$  is in star normal form. In the next section, we show that a regular expression can be transformed into star normal form, in linear time, while leaving the Glushkov automaton intact. Together, this implies that the Glushkov automaton can be constructed from an expression in quadratic time.

Let  $\Sigma$  be a finite alphabet of symbols. Uppercase letters such as  $E$ ,  $F$ , and  $G$  denote regular expressions and  $\mathcal{L}(E)$  denotes the language specified by a regular expression  $E$ . To indicate different positions or occurrences of the same symbol in an expression, we mark symbols with subscripts. For example, the regular expression  $(a + b)^* a(ab)^*$  is written as  $(a_1 + b_1)^* a_2(a_3 b_2)^*$ . With this approach the subscripted symbols  $a_i$  and  $b_j$  are called *positions* and the set of subscripted symbols in an expression  $E$  written in this form is denoted by  $\text{pos}(E)$ . Subscripting implies, for expressions  $F + G$  and  $FG$ , that  $\text{pos}(F)$  and  $\text{pos}(G)$  are disjoint.

We use  $x, y, z$  as variables for positions and  $a, b, c$  for elements of  $\Sigma$ . Finally, for a position  $x$ , let  $\chi(x)$  be the corresponding symbol of  $\Sigma$ .

The size of a regular expression  $E$  is the number of symbols it contains, including the syntactic symbols such as brackets,  $+$ ,  $\cdot$ , and  $*$ . The size of an NFA is the number of its transitions.

The following two definitions are due to Glushkov [10], who used them to define a DFA recognizing  $\mathcal{L}(E)$ . Three functions capture the notion of a position in a regular expression matching a symbol in a word. These functions are:  $\text{first}(E)$ , the set of positions that match the first symbol of some word in  $\mathcal{L}(E)$ ;  $\text{last}(E)$ , the dual set for last positions and symbols; and  $\text{follow}(E, x)$ , the set of positions that can follow position  $x$  in a path through  $E$ .

**Definition 2.1.** We can define  $first(E)$  and  $last(E)$  inductively:

$$\begin{aligned}
[E = \varepsilon \text{ or } \emptyset] \quad & first(E) = last(E) = \emptyset. \\
[E = x] \quad & first(E) = last(E) = \{x\}. \\
[E = F + G] \quad & first(E) = first(F) \cup first(G), \quad last(E) = last(F) \cup last(G). \\
[E = FG] \quad & first(E) = \begin{cases} first(F) \cup first(G) & \text{if } \varepsilon \in \mathcal{L}(F), \\ first(F) & \text{otherwise,} \end{cases} \\
& last(E) = \begin{cases} last(F) \cup last(G) & \text{if } \varepsilon \in \mathcal{L}(G), \\ last(G) & \text{otherwise.} \end{cases} \\
[E = F^*] \quad & first(E) = first(F), \quad last(E) = last(F).
\end{aligned}$$

**Definition 2.2.** The function  $follow(E, \cdot)$  maps positions of  $E$  to subsets of positions of  $E$ .

$$\begin{aligned}
[E = \varepsilon \text{ or } \emptyset] \quad & E \text{ has no positions.} \\
[E = x] \quad & follow(E, x) = \emptyset. \\
[E = F + G] \quad & follow(E, x) = \begin{cases} follow(F, x) & \text{if } x \in pos(F), \\ follow(G, x) & \text{if } x \in pos(G). \end{cases} \\
[E = FG] \quad & follow(E, x) = \begin{cases} follow(F, x) & \text{if } x \in pos(F) \setminus last(F), \\ follow(F, x) \cup first(G) & \text{if } x \in last(F), \\ follow(G, x) & \text{if } x \in pos(G). \end{cases} \\
[E = F^*] \quad & follow(E, x) = \begin{cases} follow(F, x) & \text{if } x \in pos(F) \setminus last(F), \\ follow(F, x) \cup first(F) & \text{if } x \in last(F). \end{cases}
\end{aligned}$$

Using the functions  $first$ ,  $last$ , and  $follow$ , several authors have defined the Glushkov NFA  $M_E$  recognizing  $\mathcal{L}(E)$  [5, 2, 4]. Berry and Sethi have shown that  $M_E$  is a natural representation of  $E$  [4].

**Definition 2.3.** We define the Glushkov automaton  $M_E = (Q_E \cup \{q_1\}, \Sigma, \delta_E, q_1, F_E)$  as follows:

- (1)  $Q_E = pos(E)$ , i.e. the states of  $M_E$  are the positions of  $E$  plus a new, initial state,  $q_1$ .
- (2) For  $a \in \Sigma$ , let  $\delta_E(q_1, a) = \{x \mid x \in first(E), \chi(x) = a\}$ .
- (3) For  $x \in pos(E)$ ,  $a \in \Sigma$ , let  $\delta_E(x, a) = \{y \mid y \in follow(E, x), \chi(y) = a\}$ .
- (4)  $F_E = \begin{cases} last(E) \cup \{q_1\} & \text{if } \varepsilon \in \mathcal{L}(M_E), \\ last(E) & \text{otherwise.} \end{cases}$

**Proposition 2.4.**  $\mathcal{L}(M_E) = \mathcal{L}(E)$ .

The inductive definition suggests a computation of *first*, *last*, and *follow* that is cubic in the size of  $E$ . First, we describe this canonical method. Then we refine the method to achieve quadratic time complexity.

Let  $n$  be the size of  $E$ . We begin by converting  $E$  into a syntax tree. The external nodes are labeled with  $\emptyset$ ,  $\varepsilon$ , and the occurrences of symbols, and the internal nodes are labeled with one of the operators  $+$ ,  $\cdot$ , or  $*$ . Since the regular expressions are generated by an LL(1) grammar, this can be done in time  $O(n)$  [12]. Each node  $v$  of the syntax tree corresponds to a subexpression  $E_v$  of  $E$ .

At each node  $v$  of the syntax tree we provide variables

*nullable*( $v$ ): Boolean and

*first*( $v$ ), *last*( $v$ ):  $2^{\text{pos}(E)}$ .

Furthermore, for each  $x \in \text{pos}(E)$ , there is a global variable

*follow*( $x$ ):  $2^{\text{pos}(E)}$ .

The variable *nullable*( $v$ ) indicates whether the subexpression  $E_v$  corresponding to  $v$  contains the empty word, *first*( $v$ ) and *last*( $v$ ) hold the first and last positions of  $E_v$ , and *follow*( $x$ ) holds the positions of  $E$  following  $x$  in  $E$ . We perform a postorder traversal of the syntax tree and at each node  $v$ , the variables for  $v$  are computed. More precisely, at each node  $v$  the following code is executed.

**case**

$v$  is a node labeled  $\emptyset$ :

*nullable*( $v$ ) := false;

*first*( $v$ ) :=  $\emptyset$ ;

*last*( $v$ ) :=  $\emptyset$ ;

$v$  is a node labeled  $\varepsilon$ :

*nullable*( $v$ ) := true;

*first*( $v$ ) :=  $\emptyset$ ;

*last*( $v$ ) :=  $\emptyset$ ;

$v$  is a node labeled  $x$ :

*nullable*( $v$ ) := false;

*follow*( $x$ ) :=  $\emptyset$ ;

*first*( $v$ ) :=  $\{x\}$ ;

*last*( $v$ ) :=  $\{x\}$ ;

$v$  is a node labeled  $+$ :

*nullable*( $v$ ) := *nullable*(*leftchild*) or *nullable*(*rightchild*);

*first*( $v$ ) := *first*(*leftchild*)  $\cup$  *first*(*rightchild*);

(★)

*last*( $v$ ) := *last*(*leftchild*)  $\cup$  *last*(*rightchild*);

(★)

$v$  is a node labeled  $\cdot$ :

```

  nullable( $v$ ) := nullable(leftchild) and nullable(rightchild);
  for each  $x$  in last(leftchild) do
    follow( $x$ ) := follow( $x$ )  $\cup$  first(rightchild);
  if nullable(leftchild) then
    first( $v$ ) := first(leftchild)  $\cup$  first(rightchild)
  else
    first( $v$ ) := first(leftchild);
  if nullable(rightchild) then
    last( $v$ ) := last(leftchild)  $\cup$  last(rightchild)
  else
    last( $v$ ) := last(rightchild);

```

$v$  is a node labeled  $*$ :

```

  nullable( $v$ ) := true;
  for each  $x$  in last(child) do
    follow( $x$ ) := follow( $x$ )  $\cup$  first(child);
  first( $v$ ) := first(child);
  last( $v$ ) := last(child);

```

**end case;**

**Lemma 2.5.** *The following invariant holds after node  $v$  has been visited:*

- (1) *nullable( $v$ ) is true if and only if  $\varepsilon \in \mathcal{L}(E_v)$ .*
- (2) *first( $v$ ) = first( $E_v$ ), last( $v$ ) = last( $E_v$ ).*

*Furthermore, if node  $v$  has been visited but the parent of  $v$  has not, then*

- (3) *follow( $x$ ) = follow( $E_v, x$ ) for  $x \in \text{pos}(E_v)$ .*

*Especially, for the root node  $v_0$ ,*

- (1) *first( $v_0$ ) = first( $E$ ), last( $v_0$ ) = last( $E$ ).*
- (2) *follow( $x$ ) = follow( $E, x$ ) for  $x \in \text{pos}(E)$ .*

If sets are represented as ordered lists, then the union of two sets can be implemented in time linear in the size of the sets. Since all sets are at most of size  $n$ , the algorithm to compute  $\text{first}(E)$ ,  $\text{last}(E)$ , and  $\text{follow}(E, x)$ , for  $x \in \text{pos}(E)$ , takes time  $O(n^3)$ .

The first observation on the way to a better time bound is that all unions labeled  $(*)$  or  $(**)$  are disjoint. This is because  $\text{pos}(F) \cap \text{pos}(G) = \emptyset$  if  $F + G$  or  $FG$  are subexpressions of  $E$ . Only the unions labeled  $(***)$  are not necessarily disjoint. A starred subexpression  $H^*$  of  $E$  adds the elements of  $\text{first}(H)$  to  $\text{follow}(H, x)$  for  $x \in \text{last}(H)$ , but some elements of  $\text{first}(H)$  may already belong to  $\text{follow}(H, x)$ , for some  $x \in \text{last}(H)$ , as the expression  $(a^*b^*)^*$  illustrates.

Our general strategy is as follows: We only consider expressions for which all unions, including the ones of type  $(***)$ , are disjoint. Such expressions are in star normal form. Then we show that our algorithm runs in time  $O(\text{size}(M_E))$  for expressions  $E$  in star normal form. Finally, in the next section, we show why the restriction to star normal form is justified.

**Definition 2.6.** A regular expression  $E$  is in *star normal form* if, for each starred subexpression  $H^*$  of  $E$ , the *SNF-conditions*

$$\text{follow}(H, \text{last}(H)) \cap \text{first}(H) = \emptyset$$

and

$$\varepsilon \notin \mathcal{L}(H)$$

hold.

**Lemma 2.7.** *Let  $E$  be a regular expression in star normal form. Then  $M_E$  can be computed from  $E$  in time  $O(\text{size}(E) + \text{size}(M_E))$ .*

**Proof.** Let  $E$  be in star normal form. First, let us look at the unions labeled  $(*)$ . They have the general form  $X := Y \cup Z$ , where  $Y$  and  $Z$  are disjoint. Furthermore,  $Y$  and  $Z$  will never again be referred to by the program. Thus, we can represent sets as unordered lists and we can implement the union in constant time as list concatenation without copying, possibly destroying the binding of  $Y$  and  $Z$  to its values in the process.

The unions of type  $(**)$  and  $(***)$  also have the form  $X := Y \cup Z$ , where  $Y$  and  $Z$  are disjoint. In these cases,  $Z$  is referred to several times in a **for**-loop and, thus, must be preserved. Hence, we implement the union as copying the elements of  $Z$  one by one to the end of  $Y$ . The run time is proportional to the size of  $Z$ .

Finally, we have to estimate the run time of the algorithm against the size of  $M_E$ . The crucial observation is that for any subexpression  $F$  of a subexpression  $G$  of  $E$  and for any  $x \in \text{pos}(F)$ , we have

$$\text{follow}(F, x) \subseteq \text{follow}(G, x) \subseteq \text{follow}(E, x).$$

Since all unions are disjoint, the run time spent with instruction  $(**)$  or  $(***)$  in a node  $v$  and for a position  $x$  is proportional to the number of positions in  $\text{follow}(E_v, x)$  that are not present in any of the subexpressions of  $E_v$ . Thus, the total run time spent with instructions  $(**)$  and  $(***)$  is proportional to

$$\sum_{x \in \text{pos}(E)} |\text{follow}(E, x)|,$$

which is less than or equal to the number of transitions in  $M_E$ .  $\square$

### 3. Star normal form

The goal of this section is to transform a regular expression  $E$ , in linear time, into an expression  $E^*$  in star normal form such that  $M_E = M_{E^*}$ .



**Theorem 3.1.** *For each regular expression  $E$ , there is a regular expression  $E^*$  such that*

- (1)  $M_{E^*} = M_E$ ,
- (2)  $E^*$  is in star normal form,
- (3)  $E^*$  can be computed from  $E$  in linear time.

As an intermediate step, we show that a starred expression  $E^*$  can be transformed into an expression  $E^{\circ*}$  with identical Glushkov automaton, such that the SNF-conditions of Definition 2.6 are fulfilled at least at the outermost level, namely for  $E^{\circ}$ . The crucial observation is that, if we remove from  $M_E$  all “feedback” transitions leading from final states (apart from  $q_1$ ) to states that  $q_1$  is directly connected to, and if we make  $q_1$  nonfinal, then the resulting NFA is the Glushkov automaton of an expression  $E^{\circ}$  with

$$\text{follow}(E^{\circ}, \text{last}(E^{\circ})) \cap \text{first}(E^{\circ}) = \emptyset.$$

Furthermore, all “feedback” transitions deleted from  $M_E$  in  $M_{E^{\circ}}$  are reintroduced in  $M_{E^{\circ*}}$ . Thus, we have  $M_{E^{\circ*}} = M_{E^*}$ .

**Definition 3.2.** We define  $E^{\circ}$  inductively as follows:

$$[E = \emptyset \text{ or } \varepsilon] \quad E^{\circ} = \emptyset.$$

$$[E = a] \quad E^{\circ} = E.$$

$$[E = F + G] \quad E^{\circ} = F^{\circ} + G^{\circ}.$$

$$[E = FG] \quad E^{\circ} = \begin{cases} FG & \text{if } \varepsilon \notin \mathcal{L}(F), \varepsilon \notin \mathcal{L}(G). \\ F^{\circ}G & \text{if } \varepsilon \notin \mathcal{L}(F), \varepsilon \in \mathcal{L}(G). \\ FG^{\circ} & \text{if } \varepsilon \in \mathcal{L}(F), \varepsilon \notin \mathcal{L}(G). \\ F^{\circ} + G^{\circ}(!) & \text{if } \varepsilon \in \mathcal{L}(F), \varepsilon \in \mathcal{L}(G). \end{cases}$$

$$[E = F^*] \quad E^{\circ} = F^{\circ}.$$

**Lemma 3.3.**

- (1)  $\text{size}(E^{\circ}) \leq \text{size}(E)$ .
- (2)  $\varepsilon \notin \mathcal{L}(E^{\circ})$ .
- (3)  $\text{pos}(E^{\circ}) = \text{pos}(E)$ .
- (4)  $\text{first}(E^{\circ}) = \text{first}(E)$ ,  
 $\text{last}(E^{\circ}) = \text{last}(E)$ .
- (5)  $\text{follow}(E^{\circ}, x) = \text{follow}(E, x)$  for all  $x \in \text{pos}(E) \setminus \text{last}(E)$ .
- (6)  $\text{follow}(E^{\circ}, x) = \text{follow}(E, x) \setminus \text{first}(E)$ , for all  $x \in \text{last}(E)$ , especially  
 $\text{follow}(E^{\circ}, \text{last}(E^{\circ})) \cap \text{first}(E^{\circ}) = \emptyset$ .
- (7)  $\text{follow}(E^{\circ*}, x) = \text{follow}(E^*, x)$  for all  $x \in \text{pos}(E)$ .
- (8)  $M_{E^{\circ*}} = M_{E^*}$ .

**Proof.** The first four claims are straightforward inductions on  $E$ . Claims 5 and 6 are proved by induction of  $E$ . We only show the induction step for concatenation.

$[E = FG]$ : We have the following cases.

*Case 1:*  $\varepsilon \notin \mathcal{L}(F)$ ,  $\varepsilon \notin \mathcal{L}(G)$ . In this case,  $E^\circ = E$ . For any  $x \in \text{last}(E)$ , we have  $\text{follow}(E, x) \cap \text{first}(E) = \emptyset$ . This implies  $\text{follow}(E^\circ, x) = \text{follow}(E, x) = \text{follow}(E, x) \setminus \text{first}(E)$ .

*Case 2:*  $\varepsilon \notin L(F)$ ,  $\varepsilon \in \mathcal{L}(G)$ . For any  $x \in \text{last}(G)$  we also have  $\text{follow}(E, x) \cap \text{first}(E) = \emptyset$  because  $\varepsilon \notin \mathcal{L}(F)$ . The other cases follow from the induction hypothesis.

*Case 3:*  $\varepsilon \in \mathcal{L}(F)$ ,  $\varepsilon \notin \mathcal{L}(G)$ .  $\text{last}(E^\circ) = \text{last}(G^\circ)$  because  $\varepsilon \notin \mathcal{L}(G^\circ)$ .

*Case 4:*  $\varepsilon \in \mathcal{L}(F)$ ,  $\varepsilon \in \mathcal{L}(G)$ . This case follows directly from the induction hypothesis. Claims (7) and (8) follow directly from (5) and (6).  $\square$

Substituting an expression  $H^*$  with  $H^\circ$  leaves the Glushkov automaton of  $H^*$  intact. Furthermore,

$$\text{follow}(H^\circ, \text{last}(H^\circ)) \cap \text{first}(H^\circ) = \emptyset.$$

Thus, if we substitute in  $E$  each starred subexpression  $H^*$  with  $H^\circ$ , proceeding bottom up in  $E$ , we can expect to get an expression  $E^*$  in star normal form with  $M_E = M_{E^*}$ .

**Definition 3.4.**

$$\begin{aligned} [E = \emptyset, \varepsilon, \text{ or } a] \quad & E^* = E. \\ [E = F + G] \quad & E^* = F^* + G^* \\ [E = FG] \quad & E^* = F^* G^*. \\ [E = F^*] \quad & E^* = F^{\circ\circ*}. \end{aligned}$$

**Lemma 3.5.**

- (1)  $\mathcal{L}(E) = \mathcal{L}(E^*)$ .
- (2)  $\text{size}(E^*) \leq \text{size}(E)$ .
- (3)  $\text{pos}(E^*) = \text{pos}(E)$ .
- (4)  $\text{first}(E^*) = \text{first}(E)$ ,  
 $\text{last}(E^*) = \text{last}(E)$ .
- (5)  $\text{follow}(E^*, x) = \text{follow}(E, x)$  for  $x \in \text{pos}(E)$ .
- (6)  $q_1 \in F_{E^*}$  if and only if  $q_1 \in F_E$ .

Claims (4) and (5) imply the first part of Theorem 3.1, namely  $M_{E^*} = M_E$ .

Our next claim is that  $E^*$  is in star normal form. The proof is by induction on the size of  $E$ . The interesting case is the star in the induction step.

$[E = F^*]$ : We have  $E^* = F^{**}$ . The SNF-conditions hold for  $F^{**}$  (Lemma 3.3), and the induction hypothesis implies that  $F^{**}$  is in star normal form. To complete the proof, we only have to show that  $F^{**} = F^{**}$ .

**Lemma 3.6.**

- (1)  $E^{\circ\circ} = E^{\circ}$ .
- (2)  $E^{**} = E^{**}$ .
- (3)  $E^{**} = E^*$ .

**Proof.** By induction on  $E$ . We show the induction step for the star.

$[E = F^*]$ :

(1)  $E^{\circ\circ}$  is identical to  $F^{\circ\circ}$  by definition, which, in turn, is  $F^{\circ}$  by the induction hypothesis or is  $E^{\circ}$  by definition.

(2)  $E^{**} = F^{**} = F^{**}$  by definition. Applying (1) gives  $F^{**}$ , which is  $F^{**}$  by the induction hypothesis or  $E^{**}$  by definition.

(3)  $E^{**} = F^{**} = F^{**}$  by definition. Applying (2) gives  $F^{**}$ , which is  $F^{**}$  by the induction hypothesis and (1). This, in turn, is  $E^*$  by definition.  $\square$

Finally, we show that  $E^*$  can be computed from  $E$  in linear time.  $E^*$  is built up from  $H^*$  and  $H^{\circ}$  for subexpressions  $H$  of  $E$ . Thus, we compute  $H^*$  and  $H^{\circ}$  simultaneously, during a postorder traversal through the syntax tree of  $E$ . The following lemma, together with the recursive definition of  $E^*$ , makes sure that at each node only a constant amount of time is spent. This completes the proof of Theorem 3.1.

**Lemma 3.7.**

$$[E = \emptyset \text{ or } \varepsilon] \quad \emptyset^{**} = \emptyset = \varepsilon^{**}.$$

$$[E = a] \quad E^{**} = E.$$

$$[E = F + G] \quad E^{**} = F^{**} + G^{**}.$$

$$[E = FG] \quad E^{**} = \begin{cases} F^*G^* & \text{if } \varepsilon \notin \mathcal{L}(F), \varepsilon \notin \mathcal{L}(G). \\ F^{**}G^* & \text{if } \varepsilon \notin \mathcal{L}(F), \varepsilon \in \mathcal{L}(G). \\ F^*G^{**} & \text{if } \varepsilon \in \mathcal{L}(F), \varepsilon \notin \mathcal{L}(G). \\ F^{**} + G^{**} & \text{if } \varepsilon \in \mathcal{L}(F), \varepsilon \in \mathcal{L}(G). \end{cases}$$

$$[E = F^*] \quad E^{**} = F^{**}.$$

**Proof.** By induction on  $E$ . In the concatenation step one has to observe that  $\mathcal{L}(E) = \mathcal{L}(E^*)$ , and in the star step one has to apply Lemma 3.6.  $\square$

**Example 3.8.** By definition, we have

$$(a^*b^*)^{**} = (a^*b^*)^{**}.$$

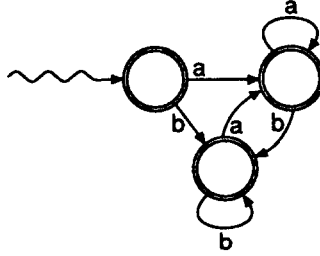


Fig. 1. The Glushkov automaton corresponding to  $(a^*b^*)^*$  and its star-normal form  $(a+b)^*$ .

Repeated application of Lemma 3.7 yields

$$\begin{aligned}
 (a^*b^*)^{\circ\circ*} &= (a^{*\circ\circ} + b^{*\circ\circ})^* \\
 &= (a^{\circ\circ} + b^{\circ\circ})^* \\
 &= (a+b)^*.
 \end{aligned}$$

Hence,  $(a+b)^*$  is the star normal form of  $(a^*b^*)^*$ . Both expressions have the same Glushkov NFA, which is shown in Fig. 1.

Putting the results of the previous two sections together, we get the following theorem.

**Theorem 3.9.** *The Glushkov automaton  $M_E$  can be computed from a regular expression  $E$  in time linear in  $\text{size}(E) + \text{size}(M_E)$ .*

**Proof.** First, we compute  $E^*$  from  $E$  in linear time. According to Theorem 3.1,  $E^*$  fulfills the precondition of Lemma 2.7. Hence, the NFA  $M_{E^*}$  can be computed from  $E^*$  in time linear in  $\text{size}(E^*) + \text{size}(M_{E^*})$ . But  $M_{E^*}$  is identical to  $M_E$ .  $\square$

Since each transition leading to a state  $x \in \text{pos}(E)$  in  $M_E$  has the label  $\chi(x)$ , the size of  $M_E$  is quadratic in the size of  $E$ . Our result is worst-case optimal, because the size of a minimal NFA equivalent to  $E$  is  $\Omega(\text{size}(E)^2)$  [15].

People writing document grammars in the SGML context are especially interested in regular expressions whose Glushkov automaton is a DFA. For such expressions, the Glushkov automaton can be constructed in linear time.

**Definition 3.10.** A regular expression  $E$  is *deterministic* if the corresponding NFA  $M_E$  is deterministic.

**Theorem 3.11.** (1) *It can be decided in linear time whether a regular expression  $E$  is deterministic.*

(2) If  $E$  is deterministic, then the deterministic finite automaton  $M_E$  can be computed from  $E$  in linear time.

**Proof.** Since the Glushkov automata of  $E$  and  $E^*$  are isomorphic,  $E$  is deterministic if and only if  $E^*$  is. Hence, we can assume that  $E$  is in star normal form. We start to compute  $first(E)$ ,  $last(E)$ , and  $follow(E, x)$  for  $x \in pos(E)$  incrementally, as outlined in the previous section, keeping track in a  $|pos(E)| \times |\Sigma|$ -matrix for which  $x \in pos(E)$  and  $a \in \Sigma$  a position  $y$  with  $\chi(y) = a$  has already been added to  $follow(E, x)$ . As soon as a position  $z$  is added to the  $follow$  set of a position  $x$  that already contains a position  $y$  with  $\chi(y) = \chi(z)$ ,  $E$  is reported as nondeterministic because  $E$  being in star normal form implies  $y \neq z$ . At this point, only time linear in the size of  $E$  has been spent.

If no such situation occurs, the entire Glushkov automaton of  $E$  is constructed. In this case,  $E$  is deterministic and the size of  $M_E$  is linear in the size of  $E$ . Thus, the time spent in constructing  $M_E$  is linear in the size of  $E$ , too.  $\square$

#### 4. Ambiguity in automata and expressions

Two types of unambiguity of regular expressions have been defined in the literature. An expression  $E$  is weakly unambiguous [5] if each word of  $E$  can be traced uniquely with a path through  $E$ , whereas  $E$  is strongly unambiguous [15] if each word of  $E$  can be uniquely decomposed into subwords according to the syntactic structure of  $E$ . The relationship between the two concepts of unambiguity has not been investigated so far. It turns out in this section that the missing link is the star normal form defined above. Thus, modulo a technical condition on the empty word, an expression  $E$  is strongly unambiguous if and only if it is weakly unambiguous and in star normal form (Theorem 4.9).

First, we define weak and strong unambiguity. From now on we consider only regular expressions that do not use  $\emptyset$  as a syntactic constituent and, hence, consider only nonempty regular languages.

**Definition 4.1.** (1) An  $\varepsilon$ -NFA  $M$  is *unambiguous* if, for each word  $w$ , there is at most one path from the initial state to a final state that spells out  $w$ .

(2) A regular expression  $E$  is *weakly unambiguous* if and only if the NFA  $M_E$  is unambiguous.

Note that a path through  $M_E$  is uniquely determined by the sequence  $x_1, \dots, x_n$  of positions in  $pos(E)$  it passes through because all transitions leading to state  $x \in pos(E)$  are labeled with  $\chi(x)$  and no transition in  $M_E$  leads to the initial state.

**Definition 4.2.** We define for languages  $L, L'$ :

(1) The concatenation of  $L$  and  $L'$  is *unambiguous* if  $v, w \in L, v', w' \in L'$ , and  $vv' = ww'$  imply  $v = w$  and  $v' = w'$ .

(2) The star of  $L$  is *unambiguous* if  $v_1, \dots, v_m \in L$ ,  $w_1, \dots, w_n \in L$ ,  $m, n \geq 0$ , and  $v_1 \dots v_m = w_1 \dots w_n$  imply  $m = n$  and  $v_i = w_i$  for  $1 \leq i \leq m$ .

**Definition 4.3.** We define inductively when a regular expression  $E$  is *strongly unambiguous*.

$[E = \varepsilon \text{ or } a]$   $E$  is strongly unambiguous.

$[E = F + G]$   $E$  is strongly unambiguous if  $F$  and  $G$  are strongly unambiguous and  $\mathcal{L}(F)$  and  $\mathcal{L}(G)$  are disjoint.

$[E = FG]$   $E$  is strongly unambiguous if  $F$  and  $G$  are strongly unambiguous and the concatenation of  $\mathcal{L}(F)$  and  $\mathcal{L}(G)$  is unambiguous.

$[E = F^*]$   $E$  is strongly unambiguous if  $F$  is strongly unambiguous and the star of  $\mathcal{L}(F)$  is unambiguous.

Strong unambiguity can be defined in terms of automata as well.

**Definition 4.4.** Let  $M'_E$  be the  $\varepsilon$ -NFA recognizing  $\mathcal{L}(E)$  according to any of the standard textbook constructions [2, 12, 14, 17, 15, 3].

**Lemma 4.5.**  $E$  is strongly unambiguous if and only if  $M'_E$  is unambiguous.

**Proof.** Sippu and Soisalon-Soininen [15] have shown this for their construction. The other variants are similar.  $\square$

**Lemma 4.6.** If  $E$  is strongly unambiguous, then  $E$  is weakly unambiguous.

**Proof.** Elimination of  $\varepsilon$ -transitions using the method of Sippu and Soisalon-Soininen [15] transforms  $M'_E$  into  $M_E$ . Thus, different paths in  $M_E$  spelling out a word  $w$  correspond to different paths in  $M'_E$  doing the same. Therefore, unambiguity of  $M'_E$  implies unambiguity of  $M_E$ . Now Lemma 4.5 can be applied.  $\square$

Now we investigate under what circumstances weakly unambiguous expressions are also strongly unambiguous. A direct comparison is facilitated through the following inductive definition of weak unambiguity.

**Lemma 4.7.**

$[E = \varepsilon \text{ or } a]$   $E$  is weakly unambiguous.

$[E = F + G]$   $E$  is weakly unambiguous if and only if  $F$  and  $G$  are weakly unambiguous and at most the empty word  $\varepsilon$  is both in  $\mathcal{L}(F)$  and  $\mathcal{L}(G)$ .

$[E = FG]$   $E$  is weakly unambiguous if and only if  $F$  and  $G$  are weakly unambiguous and the concatenation of  $\mathcal{L}(F)$  and  $\mathcal{L}(G)$  is unambiguous.

$[E = F^*]$  Let  $\text{follow}(F, \text{last}(F)) \cap \text{first}(F) = \emptyset$ ,  $\varepsilon \notin \mathcal{L}(F)$ . Then  $E$  is weakly unambiguous if and only if  $F$  is weakly unambiguous and the star of  $\mathcal{L}(F)$  is unambiguous.

**Proof.**

$[E = F + G]$  Since Glushkov automata have no  $\varepsilon$ -transitions, the only path denoting the empty word is the empty path. Furthermore, any path through  $F$  or through  $G$  is also a path through  $E$ , and any nonempty path through  $F$  is different from any path through  $G$ .

$[E = FG]$  Let us assume that  $E$  is weakly unambiguous. Since  $\mathcal{L}(F) \neq \emptyset \neq \mathcal{L}(G)$ , each path through  $F$  or  $G$  can be completed to a path through  $E$ . Thus,  $F$  and  $G$  are weakly unambiguous. Each decomposition of a word  $u \in \mathcal{L}(F)\mathcal{L}(G)$ ,  $u = vw = v'w'$  with  $v, v' \in \mathcal{L}(F)$ ,  $w, w' \in \mathcal{L}(G)$ , corresponds to paths  $x_1 \dots x_m y_1 \dots y_n$  and  $x'_1 \dots x'_{m'} y'_1 \dots y'_{n'}$  of  $E$ , where the  $x$ -positions belong to  $F$  and the  $y$ -positions to  $G$ . Since  $E$  is weakly unambiguous, the paths through  $E$  are identical. Since the positions of  $F$  and  $G$  are disjoint, we have  $m = m'$  and  $n = n'$ , i.e.  $v = v'$ ,  $w = w'$ . Thus, the concatenation of  $\mathcal{L}(F)$  and  $\mathcal{L}(G)$  is unambiguous.

This proves one direction; the other one is obvious.

$[E = F^*]$  Since  $\varepsilon \notin \mathcal{L}(E)$ , the empty word is uniquely decomposed into a sequence of words in  $\mathcal{L}(F)$ .

Any nonempty path through  $M_E$  is determined by a sequence of positions  $x_1, \dots, x_n$ ,  $n \geq 1$ , which consists of a sequence of paths through  $M_F$ .

Because  $\text{follow}(F, \text{last}(F)) \cap \text{first}(F) = \emptyset$ , the starting positions of those paths are uniquely determined. Hence, if  $E$  is weakly unambiguous, then the star of  $F$  is unambiguous.

Again, the other direction is obvious.  $\square$

Thus, weak and strong unambiguity have exactly the same inductive definition for expressions  $E$  in star normal form, provided that no subexpression of  $E$  denotes the empty word ambiguously. We call the last condition *epsilon normal form*.

**Definition 4.8.** We define by induction on  $E$  when  $E$  is in *epsilon normal form*.

$[E = \varepsilon \text{ or } a]$   $E$  is in epsilon normal form.

$[E = F + G]$   $E$  is in epsilon normal form if  $F$  and  $G$  are in epsilon normal form and  $\varepsilon \notin \mathcal{L}(F) \cap \mathcal{L}(G)$ .

$[E = FG]$   $E$  is in epsilon normal form if  $F$  and  $G$  are in epsilon normal form.

$[E = F^*]$   $E$  is in epsilon normal form if  $F$  is in epsilon normal form and  $\varepsilon \notin \mathcal{L}(F)$ .

**Theorem 4.9.**  $E$  is strongly unambiguous if and only if

- (1)  $E$  is weakly unambiguous,
- (2)  $E$  is in star normal form, and
- (3)  $E$  is in epsilon normal form.

**Proof.** Lemma 4.7 implies that for expressions in star and epsilon normal form, weak and strong unambiguity are identical. It remains to show that strongly unambiguous expressions are in star and in epsilon normal form. The crucial point in the induction is dealt with in the next lemma.  $\square$

**Lemma 4.10.** *If  $E^*$  is strongly unambiguous, then  $\text{follow}(E, \text{last}(E)) \cap \text{first}(E) = \emptyset$ .*

**Proof.** We assume that there exist  $x \in \text{last}(E)$ ,  $y \in \text{follow}(E, x) \cap \text{first}(E)$ . Since  $E$  does not contain  $\emptyset$  as a syntactic constituent, the final state  $x$  of  $M_E$  can be reached in  $M_E$  from the initial state  $q_1$  via intermediate states  $x_1, \dots, x_n$ ,  $n \geq 0$ , and some final state  $z \in \text{last}(E)$  can be reached from  $y$  via intermediate states  $y_1, \dots, y_m$ ,  $m \geq 0$ . Now the state sequence  $x_1, \dots, x_n, x, y, y_1, \dots, y_m, z$  describes a path through  $M_E$ , because  $y \in \text{follow}(E, x)$ . But this path is also the composition of two paths through  $M_E$ , because  $x \in \text{last}(E)$ ,  $y \in \text{first}(E)$ . This makes the star of  $\mathcal{L}(E)$  ambiguous.  $\square$

In the previous section, we have transformed expressions into star-normal form, in linear time. Epsilon normal form is invariant under this transformation. Thus, Theorem 4.9 reduces weak unambiguity to strong unambiguity in linear time. This yields a quadratic decision algorithm for weak unambiguity of expressions in epsilon normal form.

**Theorem 4.11.** *Regular expressions in epsilon normal form can be tested for weak unambiguity in quadratic time.*

**Proof.** Let  $E$  be in epsilon-normal form.  $E$  can be transformed into star-normal form  $E^*$  without changing the Glushkov automaton, in linear time. Furthermore,  $E^*$  is also in epsilon normal form.

Unfortunately, it is possible that  $E^*$  contains  $\emptyset$  as a syntactic constituent, even if  $E$  does not. The usual linear time elimination of  $\emptyset$  from  $E^*$ , however, resulting in an expression  $E^{*\emptyset}$ , preserves star and epsilon normal form and leaves the Glushkov automaton intact, i.e.

$$M_E = M_{E^*} = M_{E^{*\emptyset}}.$$

Now Theorem 4.9 can be applied to the  $\emptyset$ -free expression  $E^{*\emptyset}$ . Thus,  $E$  is weakly unambiguous if and only if  $E^{*\emptyset}$  is, i.e. if and only if  $E^{*\emptyset}$  is strongly unambiguous. Finally, strong unambiguity of expressions can be decided in quadratic time [15].  $\square$

## 5. Open problems

It is easy to see that a regular expression can be tested for epsilon normal form in linear time. It is an open question, however, if a given regular expression can be transformed into epsilon normal form in linear time. Our transformation into star normal form can deal with starred subexpressions. Hence, the crucial point is how expressions  $E = F + G$  with  $\varepsilon \in \mathcal{L}(F) \cap \mathcal{L}(G)$  can be handled. A straightforward approach would eliminate the empty string either from  $\mathcal{L}(F)$  or from  $\mathcal{L}(G)$ .

This opens up another question: Is there a linear-time algorithm transforming a regular expression  $E$  into an expression  $E'$  with  $\mathcal{L}(E') = \mathcal{L}(E) \setminus \{\varepsilon\}$ ?



## Acknowledgment

I thank Rolf Klein and Derick Wood for reading an earlier draft of this paper. Their thoughtful comments have considerably improved the clarity of the presentation.

## References

- [1] A.V. Aho, B.W. Kernighan and P.J. Weinberger, *The AWK Programming Language* (Addison-Wesley, Reading, MA, 1988).
- [2] A.V. Aho, R. Sethi and J.D. Ullman, *Compilers: Principles, Techniques, and Tools*, Addison-Wesley Series in Computer Science (Addison-Wesley, Reading, MA, 1986).
- [3] J. Albert and T. Ottmann, *Automaten, Sprachen und Maschinen* (Bibliographisches Institut, Mannheim, 1983).
- [4] G. Berry and R. Sethi, From regular expressions to deterministic automata, *Theoret. Comput. Sci.* **48** (1986) 117–126.
- [5] R. Book, S. Even, S. Greibach and G. Ott, Ambiguity in graphs and expressions, *IEEE Trans. Comput.* **C20** (1971) 149–153.
- [6] A. Brüggemann-Klein and D. Wood, Deterministic regular languages, in: A. Finkel and M. Jantzen, eds., Proc. *STACS '92*, Lecture Notes in Computer Science, Vol. 577 (Springer, Berlin, 1992) 173–184.
- [7] J.A. Brzozowski, Derivatives of regular expressions, *J. ACM* **11** (1964) 481–494.
- [8] C.-H. Chen and R. Paige, New theoretical and computational results for regular languages, Tech. Report 587, Courant Institute, New York University, 1992; *Third Symp. Combinatorial Pattern Matching*, accepted.
- [9] D. Dougherty and T. O'Reilly, *UNIX Text Processing* (Hayden Books, Indianapolis, 1987).
- [10] V.M. Glushkov, The abstract theory of automata, *Russian Math. Surveys*, **16** (1961) 1–53.
- [11] F.C. Hennie, *Finite-State Models for Logical Machines* (Wiley, New York, 1968).
- [12] J.E. Hopcroft and J.D. Ullman, Introduction to automata theory, languages and computation, Addison-Wesley Series in Computer Science (Addison-Wesley, Reading, MA, 1979).
- [13] ISO 8879, Information processing—text and office systems—standard generalized markup language (SGML), International Organization for Standardization, 1986.
- [14] D. Perrin, Finite automata, in: J. van Leeuwen, ed., *Handbook of Theoretical Computer Science* (Elsevier, Amsterdam and Cambridge, MA, 1990).
- [15] S. Sippu and E. Soisalon-Soininen, *Parsing Theory, Vol. 1, Languages and Parsing*, of EATCS Monographs on Theoretical Computer Science (Springer, Berlin, 1988).
- [16] G. Staubach, *UNIX-Werkzeuge zur Textmusterverarbeitung* (Springer, Berlin, 1989).
- [17] D. Wood, *Theory of Computation* (Wiley, New York, 1987).