

A SIMPLE UNIVERSAL PATTERN-MATCHING AUTOMATON

DANIEL J. BERNSTEIN

ABSTRACT. Consider an infinite non-deterministic automaton with one state \boxed{p} for each regular expression p ; transitions $\boxed{q} \xrightarrow{c} \boxed{qS}$ whenever S is a character set containing c ; and null transitions $\boxed{q} \Rightarrow \boxed{q\bar{r}}$, $\boxed{q\bar{r}r} \Rightarrow \boxed{q\bar{r}}$, $\boxed{qr} \Rightarrow \boxed{q(r' \cup r)}$, and $\boxed{qr'} \Rightarrow \boxed{q(r' \cup r)}$. If this automaton starts at the empty regular expression, then \boxed{p} recognizes exactly the language defined by p , for every p . The subautomaton affecting \boxed{p} has at most $1 + \text{len } p$ states.

1. INTRODUCTION

This paper presents a nondeterministic infinite automaton that recognizes all regular expressions simultaneously. A portion of the automaton is shown in Figure 1 below.

The automaton has one state \boxed{p} for each regular expression p , and no other states. The language recognized by \boxed{p} is exactly the language defined by p . The subautomaton affecting \boxed{p} has at most $1 + \text{len } p$ states. Here $\text{len } p$ is the number of non-parenthesis symbols in p ; for example, the length of $\bar{?}xyxz$ is 7, and the length of $((xy \cup z)z) \cup yyy$ is 9.

Is it surprising that such an automaton exists? Of course not. It is well known that, for each p , there is a nondeterministic automaton recognizing p with at most $1 + \text{len } p$ states. One can mechanically assign to each state a corresponding regular expression, and finally merge the automata for all p into a single infinite automaton that behaves as described above.

What *is* surprising about the automaton in this paper is that its definition is extremely short. There is one transition $\boxed{q} \xrightarrow{c} \boxed{qC}$ for each regular expression q , character c , and character set C containing c . There are null transitions $\boxed{q} \Rightarrow \boxed{q\bar{r}}$, $\boxed{q\bar{r}r} \Rightarrow \boxed{q\bar{r}}$, $\boxed{qr} \Rightarrow \boxed{q(r' \cup r)}$, and $\boxed{qr'} \Rightarrow \boxed{q(r' \cup r)}$, for all regular expressions q, r, r' . The automaton begins at $\boxed{()}$ where $()$ is the empty pattern. That's it.

These transitions are visibly correct, in the sense that any string recognized by \boxed{p} is in the language defined by p . It is not as easy to prove that these transitions are adequate, in the sense that any string in the language defined by p is recognized by \boxed{p} . See Theorem 4.1 and Theorem 4.2. It also takes some work to prove that the subautomaton affecting p has at most $1 + \text{len } p$ states. See Theorem 7.2. These proofs occupy the remaining sections of this paper.

Date: 20000806.

2000 Mathematics Subject Classification. Primary 68Q45.

Key words and phrases. Regular expressions, nondeterministic infinite automata.

The author was supported by the National Science Foundation under grant DMS-9600083.

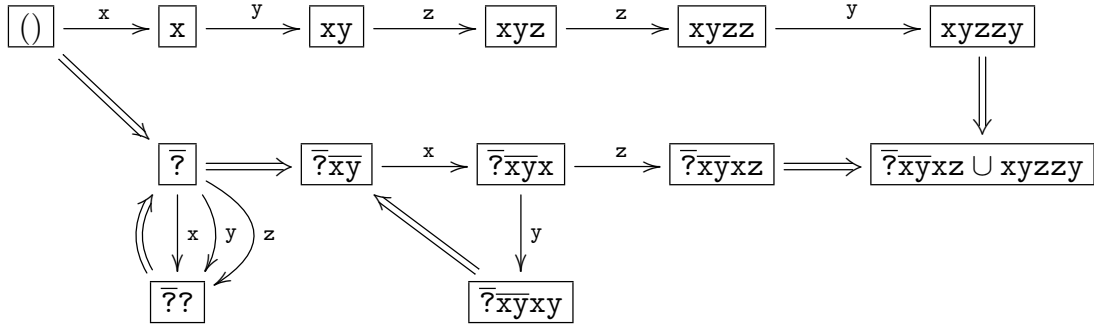


FIGURE 1. A portion of the automaton over the alphabet $\{x, y, z\}$.

Notation and terminology. The string (c_1, \dots, c_n) is abbreviated as “ $c_1 \dots c_n$ ”. For example, “ x ”, “ xyz ”, and “ $xyzzzy$ ” are strings over the alphabet $\{x, y, z\}$.

A **constant** is a set of single-character strings. The set of all single-character strings is denoted $?$. At the risk of confusion, I abbreviate the constant $\{“c”\}$ as c for each character c .

A **pattern algebra** is a set with an associative binary operation “composition” written $p, q \mapsto pq$, a neutral element for composition written $()$, a unary operation “closure” written $p \mapsto \bar{p}$, and a binary operation “union” written $p, q \mapsto p \cup q$. For example, the set of regular languages is a pattern algebra with composition $L, M \mapsto LM = \{st : s \in L, t \in M\}$; neutral element $() = \{“”\}$; union equalling the usual set union; and closure $L \mapsto \bar{L} = L^0 \cup L^1 \cup \dots$.

A **pattern** is an element of the free pattern algebra on the set of constants. In other words, a pattern is a formula built up from constants via union, closure, and composition, modulo redundant parentheses and the associativity of composition. Every pattern falls into one of the four forms $()$, qC , $q\bar{r}$, or $q(r' \cup r)$; here q , r' , and r are patterns, and C is a constant.

I write $s \in p$, and say p **matches** s , to mean that s is contained in the language defined by p . Here p is a pattern and s is a string.

Historical notes. Thompson in [1] constructed an automaton recognizing p with $O(\text{len } p)$ states. See [2] for a coherent survey of subsequent constructions. My construction is simpler than any of the constructions in [2].

I wrote down my automaton in June 1991, and distributed an implementation in a posting to `alt.sources` in January 1992. I was not familiar with the literature; the construction seemed so obvious that I assumed it was what everyone had always done. In April 1994, after reading a preliminary version of the taxonomy in [2], I announced my automaton in a posting to `comp.theory`.

2. PATTERN IMPLICATION

The relation “ $p' \Rightarrow p$ ” is the transitive closure of the **basic implications** “ $q \Rightarrow q\bar{r}$ ”, “ $q\bar{r}r \Rightarrow q\bar{r}$ ”, “ $qr \Rightarrow q(r' \cup r)$ ”, and “ $qr' \Rightarrow q(r' \cup r)$ ”. In other words, $p' \Rightarrow p$ if and only if there is a chain $p' = p_0 \Rightarrow p_1 \Rightarrow \dots \Rightarrow p_n = p$ of basic implications. In particular $p \Rightarrow p$.

For example, $() \Rightarrow \bar{?}$, and $\bar{?} \Rightarrow \bar{?x\bar{y}}$, so $() \Rightarrow \bar{?x\bar{y}}$.

Theorem 2.1. *If $p' \Rightarrow p$ and $s \in p'$ then $s \in p$.*

Proof. It suffices to check the basic implications: if $s \in q$ or $s \in q\bar{r}r$ then $s \in q\bar{r}$; and if $s \in qr$ or $s \in qr'$ then $s \in q(r' \cup r)$. \square

Theorem 2.2. *If $p' \Rightarrow p$ then $p''p' \Rightarrow p''p$.*

Proof. This is a purely formal consequence of the definition: $p''q\bar{r}r \Rightarrow p''q\bar{r}$, $p''q \Rightarrow p''q\bar{r}$, $p''qr \Rightarrow p''q(r' \cup r)$, and $p''qr' \Rightarrow p''q(r' \cup r)$, so given a chain of basic implications we may prepend p'' to each term. \square

3. THE impl FUNCTION

Define recursively

$$\text{impl } p = \begin{cases} \{p\} & \text{if } p = () \text{ or } p = qC, \\ \text{impl } qr' \cup \text{impl } qr & \text{if } p = q(r' \cup r), \\ \text{impl } q \cup \{q\bar{r}y : y \in \text{impl } r, y \neq ()\} & \text{if } p = q\bar{r}. \end{cases}$$

For example, $\text{impl } \bar{?x}\bar{y} = \{(), \bar{?}, \bar{?x}\bar{y}xy\}$.

Theorem 3.1. *If $x \in \text{impl } p$ then $x \Rightarrow p$.*

Proof. Induct on $\text{len } p$. First, if $p = ()$ or $p = qC$ then $x = p$ so $x \Rightarrow p$.

Second, say $p = q(r' \cup r)$. If $x \in \text{impl } qr'$ then by induction $x \Rightarrow qr' \Rightarrow p$. If $x \in \text{impl } qr$ then by induction $x \Rightarrow qr \Rightarrow p$.

Third, say $p = q\bar{r}$. If $x \in \text{impl } q$ then by induction $x \Rightarrow q \Rightarrow p$. If $x = q\bar{r}y$, with $y \in \text{impl } r$, then by induction $y \Rightarrow r$, so $x = q\bar{r}y \Rightarrow q\bar{r}r \Rightarrow p$ by Theorem 2.2. \square

Theorem 3.2. *If $x \in \text{impl } p$ then x is empty or ends with a constant.*

Proof. Induct on $\text{len } p$. First, if $p = ()$ or $p = qC$ then $x = p$.

Second, if $p = q(r' \cup r)$ then $x \in \text{impl } qr'$ or $x \in \text{impl } qr$; by induction x is empty or ends with a constant.

Third, say $p = q\bar{r}$. If $x \in \text{impl } q$ then by induction x is empty or ends with a constant. If $x = q\bar{r}y$, with $() \neq y \in \text{impl } r$, then by induction y ends with a constant, so x does too. \square

Theorem 3.3. *If $s \in p$ then $s \in x$ for some $x \in \text{impl } p$.*

Proof. Induct on $\text{len } p$. First, if $p = ()$ or $p = qC$ then $p \in \text{impl } p$ so there is nothing to prove.

Second, if $p = q(r' \cup r)$ then $s \in qr'$ or $s \in qr$. So by induction $s \in x$ where $x \in \text{impl } qr'$ or $x \in \text{impl } qr$. Either way $x \in \text{impl } p$.

Third, if $p = q\bar{r}$ then $s = t't$ with $t' \in q$, $t \in \bar{r}$. If $t = \text{""}$ then $s = t' \in q$ so by induction there is an $x \in \text{impl } q$ with $s \in x$; and $x \in \text{impl } p$. If $t \neq \text{""}$ we may write $t = uv$ where $u \in \bar{r}$ and $v \in r$, $v \neq \text{""}$. By induction there is a $y \in \text{impl } r$ with $v \in y$; since v is nonempty we cannot have $y = ()$. Thus $q\bar{r}y \in \text{impl } p$, and $s = t'uv \in q\bar{r}y$ as desired. \square

4. CONSEQUENCES OF impl

Theorem 4.1 says that my automaton works for the empty string. Theorem 4.2 says that, if my automaton works for a string t , then it also works for t plus any character.

Theorem 4.1. *"" $\in p$ if and only if $() \Rightarrow p$.*

Proof. Say $() \Rightarrow p$. Certainly $"" \in ()$; by Theorem 2.1, $"" \in p$.

Conversely, say $"" \in p$. By Theorem 3.3, $"" \in x$ for some $x \in \text{impl } p$. Since $"" \in x$, x cannot end with a constant; thus, by Theorem 3.2, $x = ()$. By Theorem 3.1, $() = x \Rightarrow p$. \square

Theorem 4.2. $t“c” \in p$ if and only if there exist q and C such that $t \in q$, $“c” \in C$, and $qC \Rightarrow p$.

Proof. If $t \in q$ and $“c” \in C$ then $t“c” \in qC$; if also $qC \Rightarrow p$ then $t“c” \in p$ by Theorem 2.1.

Conversely, say $t“c” \in p$. By Theorem 3.3, $t“c” \in x$ for some $x \in \text{impl } p$. Since $t“c” \neq ""$, x cannot be empty; thus, by Theorem 3.2, x ends with a constant. Write $x = qC$. Then $t \in q$ and $“c” \in C$. Finally $qC = x \Rightarrow p$ by Theorem 3.1. \square

5. LEFT SUBPATTERNS

The relation “ p' is a left subpattern of p ” is the transitive closure of the relations “ q is a left subpattern of qC ”, “ $q\bar{r}r$ is a left subpattern of $q\bar{r}$ ”, “ q is a left subpattern of $q\bar{r}$ ”, “ qr is a left subpattern of $q(r' \cup r)$ ”, and “ qr' is a left subpattern of $q(r' \cup r)$ ”.

In other words, p' is a left subpattern of p if and only if $\boxed{p'}$ affects \boxed{p} in my automaton. Therefore, if $s \in p$, then any prefix s' of s satisfies $s' \in p'$ for some left subpattern p' of p .

Example: $\bar{?}?$ is a left subpattern of $\bar{?}$, which is a left subpattern of $\bar{?}\bar{x}\bar{y}$, which is a left subpattern of $\bar{?}\bar{x}\bar{y}x$, which is a left subpattern of $\bar{?}\bar{x}\bar{y}xz$. So $\bar{?}?$ is a left subpattern of $\bar{?}\bar{x}\bar{y}xz$.

6. THE left FUNCTION

Define $\text{left}_0 p = p$. Define $\text{left}_{n+1} p$ for each n as follows:

$$\text{left}_{n+1} p = \begin{cases} \text{undefined} & \text{if } p = () \\ \text{left}_n q & \text{if } p = qC \\ \text{left}_n q\bar{r}r & \text{if } p = q\bar{r}, n < \text{len } r \\ \text{left}_{n-\text{len } r} q & \text{if } p = q\bar{r}, n \geq \text{len } r \\ \text{left}_n qr & \text{if } p = q(r' \cup r), n < \text{len } r \\ \text{left}_{n-\text{len } r} qr' & \text{if } p = q(r' \cup r), n \geq \text{len } r. \end{cases}$$

For example, the values of $\text{left}_n \bar{?}\bar{x}\bar{y}xz$, as n increases from 0 through 7, are $\bar{?}\bar{x}\bar{y}xz$, $\bar{?}\bar{x}\bar{y}x$, $\bar{?}\bar{x}\bar{y}$, $\bar{?}\bar{x}\bar{y}xy$, $\bar{?}\bar{x}\bar{y}x$, $\bar{?}$, $\bar{?}?$, and $()$. For $n > 7$, $\text{left}_n \bar{?}\bar{x}\bar{y}xz$ is undefined.

The reader may enjoy verifying certain facts not needed in this paper: $\text{left}_{\text{len } p} p = ()$; if $n < \text{len } p$ then $\text{left}_n p$ is defined and nonempty; if $n \leq \text{len } p$ then $\text{left}_n p' = p' \text{left}_n p$; if $\text{left}_n p$ is defined then it is a left subpattern of p .

Theorem 6.1. *If $n > \text{len } p$ then $\text{left}_n p$ is undefined.*

Proof. Induct on n . Note that $n > 0$.

If $p = ()$ then $\text{left}_n p$ is undefined. If $p = qC$ then $n > n - 1 > \text{len } q$ so $\text{left}_{n-1} q$ is undefined by induction; so $\text{left}_n p = \text{left}_{n-1} q$ is undefined. If $p = q\bar{r}$ then $n > n - 1 - \text{len } r > \text{len } q$ so $\text{left}_{n-1-\text{len } r} q$ is undefined by induction; so $\text{left}_n p = \text{left}_{n-1-\text{len } r} q$ is undefined. If $p = q(r' \cup r)$ then $n > n - 1 - \text{len } r > \text{len } qr'$ so $\text{left}_{n-1-\text{len } r} qr'$ is undefined by induction; so $\text{left}_n p = \text{left}_{n-1-\text{len } r} qr'$ is undefined. \square

Theorem 6.2. $\text{left}_n p'p = \text{left}_{n-\text{len } p} p'$ if $n \geq \text{len } p$.

Proof. Induct on n .

If $p = ()$ then $\text{len } p = 0$ and $p'p = p'$ as desired. If $p = qC$ then $n > n - 1 \geq \text{len } q$ so

$$\text{left}_n p'p = \text{left}_{n-1} p'q = \text{left}_{n-1-\text{len } q} p' = \text{left}_{n-\text{len } p} p'$$

by induction. If $p = q\bar{r}$ then $n > n - 1 - \text{len } r \geq \text{len } q$ so

$$\text{left}_n p'p = \text{left}_{n-1-\text{len } r} p'q = \text{left}_{n-1-\text{len } r-\text{len } q} p' = \text{left}_{n-\text{len } p} p'$$

by induction. If $p = q(r' \cup r)$ then $n > n - 1 - \text{len } r \geq \text{len } qr'$ so

$$\text{left}_n p'p = \text{left}_{n-1-\text{len } r} p'qr' = \text{left}_{n-1-\text{len } r-\text{len } qr'} p' = \text{left}_{n-\text{len } p} p'$$

by induction. □

7. THE J FUNCTION

Theorem 7.2 states that all left subpatterns of p are values of $\text{left}_n p$ for $n \in \{0, 1, \dots, \text{len } p\}$. Thus the number of states affecting \boxed{p} in my automaton is at most $1 + \text{len } p$.

Define $J(p, m, 0) = m$. For each n define $J(p, m, n + 1) =$

$$\left\{ \begin{array}{ll} \text{undefined} & \text{if } p = () \\ J(q, m, n) + 1 & \text{if } p = qC \\ J(q\bar{r}r, m, n) + 1 & \text{if } p = q\bar{r}, n < \text{len } r, J(q\bar{r}r, m, n) < \text{len } r \\ J(q\bar{r}r, m, n) - \text{len } r & \text{if } p = q\bar{r}, n < \text{len } r, J(q\bar{r}r, m, n) \geq \text{len } r \\ J(q, m, n - \text{len } r) + \text{len } r + 1 & \text{if } p = q\bar{r}, n \geq \text{len } r \\ J(qr, m, n) + 1 & \text{if } p = q(r' \cup r), n < \text{len } r, J(qr, m, n) < \text{len } r \\ J(qr, m, n) + \text{len } r' + 1 & \text{if } p = q(r' \cup r), n < \text{len } r, J(qr, m, n) \geq \text{len } r \\ J(qr', m, n - \text{len } r) + \text{len } r + 1 & \text{if } p = q(r' \cup r), n \geq \text{len } r. \end{array} \right.$$

Theorem 7.1. $\text{left}_m \text{left}_n p = \text{left}_{J(p, m, n)} p$ if $\text{left}_n p$ is defined.

Proof. If $n = 0$ then $J(p, m, n) = m$ and $\text{left}_n p = p$; thus $\text{left}_{J(p, m, n)} p = \text{left}_m p = \text{left}_m \text{left}_n p$.

Now induct on n . Assume $\text{left}_{n+1} p$ is defined. Note that $n + 1 \leq \text{len } p$ by Theorem 6.1, so $p \neq ()$. I will show that $\text{left}_m \text{left}_{n+1} p = \text{left}_{J(p, m, n+1)} p$.

1. Say $p = qC$. Write $j = J(q, m, n)$. Then

$$\text{left}_m \text{left}_{n+1} p = \text{left}_m \text{left}_n q = \text{left}_j q = \text{left}_{j+1} p.$$

2. Say $p = q\bar{r}$ and $n < \text{len } r$. Write $j = J(q\bar{r}r, m, n)$. Observe that

$$\text{left}_m \text{left}_{n+1} p = \text{left}_m \text{left}_n q\bar{r}r = \text{left}_j q\bar{r}r.$$

If $j < \text{len } r$ then $\text{left}_j q\bar{r}r = \text{left}_{j+1} p$; else $\text{left}_j q\bar{r}r = \text{left}_{j-\text{len } r} q\bar{r} = \text{left}_{j-\text{len } r} p$.

3. Say $p = q\bar{r}$ and $n \geq \text{len } r$. Write $j = J(q, m, n - \text{len } r)$. Then

$$\text{left}_m \text{left}_{n+1} p = \text{left}_m \text{left}_{n-\text{len } r} q = \text{left}_j q = \text{left}_{j+\text{len } r+1} p.$$

4. Say $p = q(r' \cup r)$ and $n < \text{len } r$. Write $j = J(qr, m, n)$. Observe that

$$\text{left}_m \text{left}_{n+1} p = \text{left}_m \text{left}_n qr = \text{left}_j qr.$$

If $j < \text{len } r$ then $\text{left}_j qr = \text{left}_{j+1} p$; otherwise

$$\text{left}_j qr = \text{left}_{j-\text{len } r} q = \text{left}_{j-\text{len } r+\text{len } r'} qr' = \text{left}_{j+\text{len } r'+1} p.$$

5. Say $p = q(r' \cup r)$ and $n \geq \text{len } r$. Write $j = J(qr', m, n - \text{len } r)$. Then $\text{left}_m \text{left}_{n+1} p = \text{left}_m \text{left}_{n-\text{len } r} qr' = \text{left}_j qr' = \text{left}_{j+\text{len } r+1} p$.

□

Theorem 7.2. *If p' is a left subpattern of p then $p' = \text{left}_n p$ for some n .*

Proof. If $p = qC$ then $q = \text{left}_1 p$.

If $p = q\bar{r}$ then $q = \text{left}_{\text{len } r+1} p$. Furthermore $q\bar{r}r = \text{left}_1 p$ if $r \neq ()$; $q\bar{r}r = \text{left}_0 p$ if $r = ()$.

If $p = q(r' \cup r)$ then $qr' = \text{left}_{\text{len } r+1} p$. Furthermore $qr = \text{left}_1 p$ if $r \neq ()$. If $r = ()$ then $qr = q = \text{left}_{\text{len } r'} qr' = \text{left}_{\text{len } r'+1} p$ by Theorem 6.2.

Finally, by Theorem 7.1, the relation “ p' is $\text{left}_n p$ for some n ” is transitive. □

REFERENCES

- [1] Ken Thompson, *Regular expression search algorithm*, Communications of the ACM **11** (1968), 419–422.
- [2] Bruce W. Watson, *Taxonomies and toolkits of regular language algorithms*, Ph.D. thesis, Eindhoven University of Technology, 1995.

DEPARTMENT OF MATHEMATICS, STATISTICS, AND COMPUTER SCIENCE, THE UNIVERSITY OF ILLINOIS AT CHICAGO, CHICAGO, IL 60607–7045

E-mail address: `djb@pobox.com`