

## DEFINITION OF REGULAR EXPRESSIONS

Using concatenations, alternation and Kleene closure, we can define *regular expressions* as follows:

- $\emptyset$  is a regular expression denoting the empty set (the set containing no strings).  $\emptyset$  is rarely used, but is included for completeness.
- $\lambda$  is a regular expression denoting the set that contains only the empty string. This set is not the same as the empty set, because it contains one element.
- A string  $s$  is a regular expression denoting a set containing the single string  $s$ .

- If  $A$  and  $B$  are regular expressions, then  $A \mid B$ ,  $AB$ , and  $A^*$  are also regular expressions, denoting the alternation, concatenation, and Kleene closure of the corresponding regular sets.

Each regular expression denotes a set of strings (a *regular set*). Any finite set of strings can be represented by a regular expression of the form  $(s_1 \mid s_2 \mid \dots \mid s_k)$ . Thus the reserved words of ANSI C can be defined as  $(\text{auto} \mid \text{break} \mid \text{case} \mid \dots)$ .

The following additional operations useful. They are not strictly necessary, because their effect can be obtained using alternation, concatenation, Kleene closure:

- $P^+$  denotes all strings consisting of *one* or more strings in  $P$  concatenated together:  
 $P^+ = (P^+ \mid \lambda)$  and  $P^+ = P P^*$ .  
For example,  $(0 \mid 1)^+$  is the set of all strings containing one or more bits.
- If  $A$  is a set of characters,  $\text{Not}(A)$  denotes  $(\Sigma - A)$ ; that is, all *characters* in  $\Sigma$  *not* included in  $A$ . Since  $\text{Not}(A)$  can never be larger than  $\Sigma$  and  $\Sigma$  is finite,  $\text{Not}(A)$  must also be finite, and is therefore regular.  $\text{Not}(A)$  does not contain  $\lambda$  since  $\lambda$  is not a character (it is a zero-length string).

For example,  $\text{Not}(\text{Eol})$  is the set of all characters excluding  $\text{Eol}$  (the end of line character,  $\backslash n$  in Java or C).

- It is possible to extend  $\text{Not}$  to strings, rather than just  $\Sigma$ . That is, if  $S$  is a set of strings, we define  $\bar{S}$  to be  $(\Sigma^* - S)$ ; the set of all strings except those in  $S$ . Though  $\bar{S}$  is usually infinite, it is also regular if  $S$  is.
- If  $k$  is a constant, the set  $A^k$  represents all strings formed by concatenating  $k$  (possibly different) strings from  $A$ . That is,  $A^k = (A A A \dots)$  ( $k$  copies). Thus  $(0 \mid 1)^{32}$  is the set of all bit strings exactly 32 bits long.

## Examples

Let  $D$  be the ten single digits and let  $L$  be the set of all 52 letters. Then

- A Java or C++ single-line comment that begins with `//` and ends with `Eol` can be defined as:

$\text{Comment} = // \text{Not}(\text{Eol})^* \text{Eol}$

- A fixed decimal literal (e.g., `12.345`) can be defined as:

$\text{Lit} = D^+ . D^+$

- An optionally signed integer literal can be defined as:

$\text{IntLiteral} = ( '+' | '-' | \lambda ) D^+$

(Why the quotes on the plus?)

- A comment delimited by `##` markers, which allows single `#`'s within the comment body:

$\text{Comment2} =$

$## ((\# | \lambda) \text{Not}(\#))^* ##$

All finite sets and many infinite sets are regular. But not all infinite sets are regular. Consider the set of balanced brackets of the form  $[[[...]]]$ .

This set is defined formally as

$\{ [^m ]^m \mid m \geq 1 \}$ .

This set is known *not* to be regular. Any regular expression that tries to define it either does not get *all* balanced nestings or it includes extra, unwanted strings.

## FINITE AUTOMATA AND SCANNERS

A *finite automaton* (FA) can be used to recognize the tokens specified by a regular expression. FAs are simple, idealized computers that recognize strings belonging to regular sets. An FA consists of:

- A finite set of *states*
- A set of *transitions* (or *moves*) from one state to another, labeled with characters in  $\Sigma$
- A special state called the *start* state
- A subset of the states called the *accepting*, or *final*, states

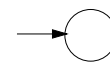
These four components of a finite automaton are often represented graphically:



is a state



is a transition



is the start state

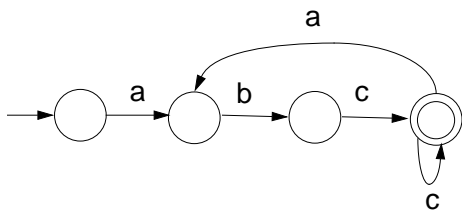


is an accepting state

Finite automata (the plural of automaton is automata) are represented graphically using *transition diagrams*. We start at the start state. If the next input character matches the label on

a transition from the current state, we go to the state it points to. If no move is possible, we stop. If we finish in an accepting state, the sequence of characters read forms a *valid* token; otherwise, we have not seen a valid token.

In this diagram, the valid tokens are the strings described by the regular expression  $(a b (c)^+)^+$ .



## DETERMINISTIC FINITE AUTOMATA

As an abbreviation, a transition may be labeled with more than one character (for example,  $\text{Not}(c)$ ). The transition may be taken if the current input character matches any of the characters labeling the transition.

If an FA always has a *unique* transition (for a given state and character), the FA is *deterministic* (that is, a deterministic FA, or DFA). Deterministic finite automata are easy to program and often drive a scanner.

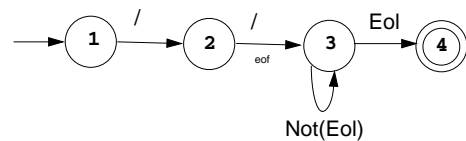
If there are transitions to more than one state for some character, then the FA is *nondeterministic* (that is, an NFA).

A DFA is conveniently represented in a computer by a *transition table*. A transition table,  $T$ , is a two dimensional array indexed by a DFA state and a vocabulary symbol.

Table entries are either a DFA state or an error flag (often represented as a blank table entry). If we are in state  $s$ , and read character  $c$ , then  $T[s,c]$  will be the next state we visit, or  $T[s,c]$  will contain an error marker indicating that  $c$  cannot extend the current token. For example, the regular expression

`// Not(Eol)* Eol`

which defines a Java or C++ single-line comment, might be translated into



The corresponding transition table is:

State	Character				
	/	Eol	a	b	...
1	2				
2	3				
3	3	4	3	3	3
4					

A complete transition table contains one column for each character. To save space, *table compression* may be used. Only non-error entries are explicitly represented in the table, using hashing, indirection or linked structures.

All regular expressions can be translated into DFAs that accept (as valid tokens) the strings defined by the regular expressions. This translation can be done manually by a programmer or automatically using a scanner generator.

A DFA can be coded in:

- Table-driven form
- Explicit control form

In the table-driven form, the transition table that defines a DFA's actions is explicitly represented in a run-time table that is "interpreted" by a driver program.

In the direct control form, the transition table that defines a DFA's actions appears implicitly as the control logic of the program.

For example, suppose **CurrentChar** is the current input character. End of file is represented by a special character value, **eof**. Using the DFA for the Java comments shown earlier, a table-driven scanner is:

```
State = StartState
while (true){
    if (CurrentChar == eof)
        break
    NextState =
        T[State][CurrentChar]
    if (NextState == error)
        break
    State = NextState
    read(CurrentChar)
}
if (State in AcceptingStates)
    // Process valid token
else // Signal a lexical error
```

This form of scanner is produced by a scanner generator; it is definition-independent. The scanner is a driver that can scan *any* token if T contains the appropriate transition table.

Here is an explicit-control scanner for the same comment definition:

```
if (CurrentChar == '/') {
    read(CurrentChar)
    if (CurrentChar == '/')
        repeat
            read(CurrentChar)
        until (CurrentChar in
            {eol, eof})
    else //Signal lexical error
else // Signal lexical error
if (CurrentChar == eol)
    // Process valid token
else //Signal lexical error
```

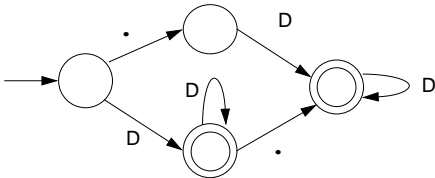
The token being scanned is "hardwired" into the logic of the code. The scanner is usually easy to read and often is more efficient, but is specific to a single token definition.

## MORE EXAMPLES

- A FORTRAN-like real literal (which requires digits on either or both sides of a decimal point, or just a string of digits) can be defined as

$$\text{RealLit} = (D^+ (\lambda | \cdot)) | (D^* \cdot D^+)$$

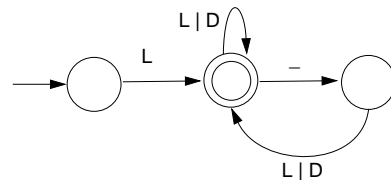
This corresponds to the DFA



- An identifier consisting of letters, digits, and underscores, which begins with a letter and allows no adjacent or trailing underscores, may be defined as

$$\text{ID} = L (L | D)^* ( \_ (L | D)^+ )^*$$

This definition includes identifiers like `sum` or `unit_cost`, but excludes `_one` and `two_` and `grand__total`. The DFA is:



## Lex/Flex/JLex

Lex is a well-known Unix scanner generator. It builds a scanner, in C, from a set of regular expressions that define the tokens to be scanned.

Flex is a newer and faster version of Lex.

JLex is a Java version of Lex. It generates a scanner coded in Java, though its regular expression definitions are very close to those used by Lex and Flex.

Lex, Flex and JLex are largely *non-procedural*. You don't need to tell the tools *how* to scan. All you need to tell it *what* you want scanned (by giving it definitions of valid tokens).

This approach greatly simplifies building a scanner, since most of the details of scanning (I/O, buffering, character matching, etc.) are automatically handled.