

## Unit-2

2.1	LEXICAL ANALYZER.....	1
2.1.1	Role of Lexical Analysis.....	1
2.1.2	Tokens, Lexemes and Patterns.....	2
2.1.3	Input Buffering.....	3
2.1.4	Specification of Tokens .....	4
2.1.5	Recognition of Tokens.....	7
2.1.6	Finite Automata.....	7
2.1.7	Design of Lexical Analyzer .....	9
2.1.8	State Minimization in DFA.....	12

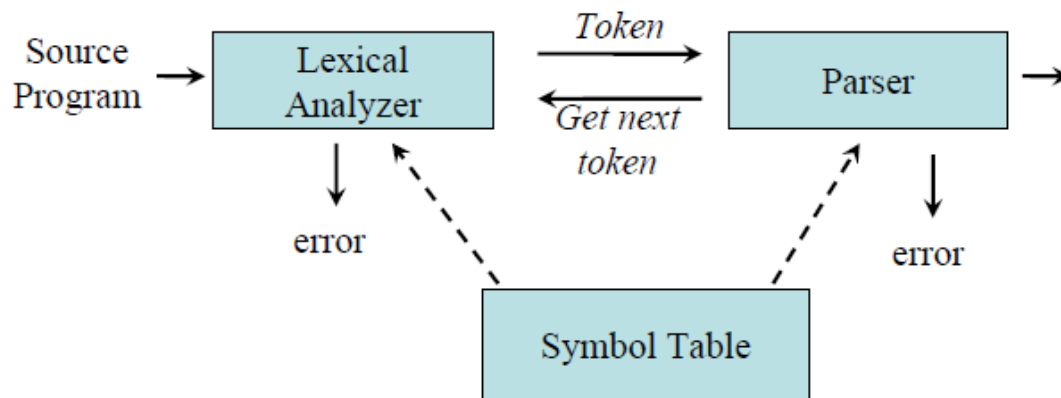
## Unit-2

### 2.1 Lexical Analyzer

#### 2.1.1 Role of Lexical Analysis

As the first phase of a compiler, the main task of the lexical analyzer is to read the input characters of the source program, group them into lexemes, and produce as output a sequence of tokens for each lexeme in the source program. The stream of tokens is sent to the parser for syntax analysis. It is common for the lexical analyzer to interact with the symbol table as well. When the lexical analyzer discovers a lexeme constituting an identifier, it needs to enter that lexeme into the symbol table. In some cases, information regarding the kind of identifier may be read from the symbol table by the lexical analyzer to assist it in determining the proper token it must pass to the parser. Another main task of the lexical analyzer is to strip out white spaces and other unnecessary character from the source program while generating the tokens. Last but not the least another main task is to correlate error message from the compiler with the source program.

These interactions are suggested in Figure. Commonly, the interaction is implemented by having the parser call the lexical analyzer. The call, suggested by the getNextToken command, causes the lexical analyzer to read characters from its input until it can identify the next lexeme and produce for it the next token, which it returns to the parser.



Since the lexical analyzer is the part of the compiler that reads the source text, it may perform certain other tasks besides identification of lexemes. One such task is stripping out comments and whitespace (blank, newline, tab, and perhaps other characters that are used to separate tokens in the input). Another task is correlating error messages generated by the compiler with the source

program. For instance, the lexical analyzer may keep track of the number of newline characters seen, so it can associate a line number with each error message. In some compilers, the lexical analyzer makes a copy of the source program with the error messages inserted at the appropriate positions. If the source program uses a macro – preprocessor, the expansion of macros may also be performed by the lexical analyzer.

### **Why separate Lexical Analyzer is needed?**

- a) Simplicity of design is the most important consideration. The separation of lexical and syntactic analysis often allows us to simplify at least one of these tasks.
- b) Compiler efficiency is improved. A separate lexical analyzer allows us to apply specialized techniques that serve only the lexical task, not the job of parsing.
- c) Specialized buffering techniques for reading input characters can speed up the compiler significantly.
- d) Compiler portability is enhanced. Input – Device – Specific peculiarities can be restricted to the lexical analyzer.

### **2.1.2 Tokens, Lexemes and Patterns**

A token is a pair consisting of a token name and an optional attribute value. The token name is an abstract symbol representing a kind of lexical unit, e.g., a particular keyword, or a sequence of input characters denoting an identifier. The token names are the input symbols that the parser processes. In what follows, we shall generally write the name of a token in boldface. We will often refer to a token by its token name.

A pattern is a description of the form that the lexemes of a token may take. In the case of a keyword as a token, the pattern is just the sequence of characters that form the keyword. For identifiers and some other tokens, the pattern is a more complex structure that is matched by many strings.

A lexeme is a sequence of characters in the source program that matches the pattern for a token and is identified by the lexical analyzer as an instance of that token.

#### **2.1.2.1 Attribute of Tokens**

When more than one lexeme can match a pattern, the lexical analyzer must provide the subsequent compiler phases additional information about the particular lexeme that matched. For example, the pattern for token number matches both 0 and 1, but it is extremely important for the code generator to know which lexeme was found in the source program. Thus, in many cases the lexical analyzer

returns to the parser not only a token name, but an attribute value that describes the lexeme represented by the token; the token name influences parsing decisions, while the attribute value influences translation of tokens after the parse.

The token name and associated attribute values for the following statement.

A=B\*C+2

<id, pointer to symbol table entry for A>

<assignment\_op>

<id, pointer to symbol table entry for B>

<mult\_op>

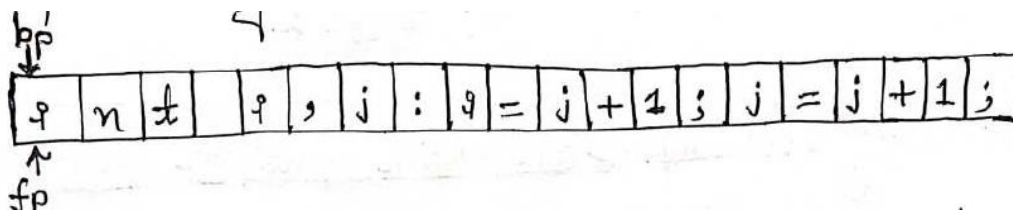
<id, pointer to symbol-table entry for C>

<add\_op>

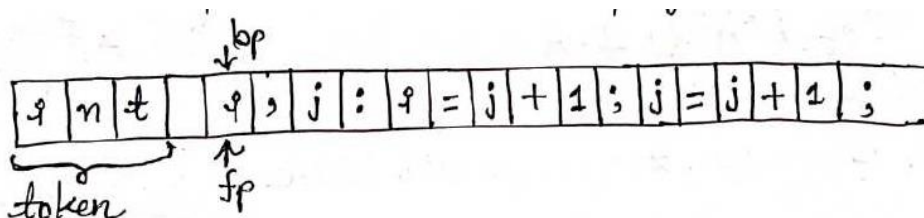
<number, integer value 2>

### 2.1.3 Input Buffering

Reading character by character from secondary storage is slow process and time consuming, so we use buffer technique to eliminate this problem and increase efficiency. The lexical analyzer scans the input from left to right one character at a time. It uses two pointers” bp” (begin pointer) and” fp” (forward pointer) to keep track of the pointer of input scanned. Initially both the pointers point to the first character of the input string.



The forward pointer moves ahead to search for end of lexeme. As soon as the blank space is encountered, it indicates end of the lexeme. In above example as soon as fp encounters a blank space the lexeme 'int' is identified. Then both bp and fp are set at next token as shown in the figure below. This process will be repeated for the whole program.



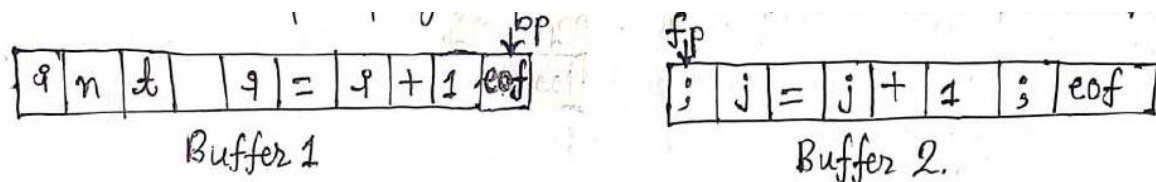
### 2.1.3.1 One Buffer Scheme:

In this scheme, only one buffer is used to store the input string. The problem with this scheme is that if lexeme is very long then it crosses the buffer boundary, to scan the rest of the lexeme the buffer must be refilled.

### 2.1.3.2 Two Buffer Scheme:

In this scheme two buffer are used to store input string and they are scanned alternately. When end of the current buffer is reached the other buffer is filled. The only problem with this method is that if length of the lexeme is longer than length of the buffer then scanning input cannot be scanned completely.

Initially both the bp and fp are pointing to the first character of first buffer. Then the fp move forward in search of end of lexeme. As soon as blank character is recognized, the string between bp and fp is identified as corresponding token. To identify, the boundary of first buffer end of buffer character (eof) should be placed at the end of the first buffer, similarly for second buffer also. Alternatively, both the buffers can be filled up until end of the input program and stream of tokens is identified.



### 2.1.4 Specification of Tokens

#### 2.1.4.1 Alphabet:

An alphabet is a finite, nonempty set of symbol. Conventionally, we use the symbol  $\Sigma$  for an alphabet. Common alphabet includes:

$\Sigma = \{0,1\}$ , this is the binary alphabet.

$\Sigma = \{a,b,\dots,z\}$ , this is the set of all lower-case letters.

#### Power of an alphabet

If  $\Sigma$  is an alphabet, we can express the set of all strings of a certain length from that alphabet using an exponential notation. We define  $\Sigma^k$  to be the set of strings of length  $k$ , each of whose symbol is in  $\Sigma$

For example: if  $\{a,b,c\}$

$\Sigma^0 = \epsilon$

$$\Sigma^1 = \{a, b, c\}$$

The set of all string over an alphabet is conventionally denoted by  $\Sigma^*$

Example  $\Sigma = \{0, 1\}$

$$\Sigma^* = \{0, 1\}^* = \{\epsilon, 0, 1, 00, 01, 10, 11, \dots\}$$

Sometimes, we wish to exclude the empty string from the set of strings. The set of nonempty string from alphabet  $\Sigma$  is denoted by  $\Sigma^+$ .

### Concentration of String

Let  $x$  and  $y$  be strings. Then  $xy$  denotes the concentration of  $x$  and  $y$ , that is the string formed by making a copy of  $x$  and following it by a copy of  $y$ .

Example: let  $x=1101$  and  $y=0011$  then  $xy=11010011$

#### 2.1.4.2 Strings:

A string is a finite sequence of symbols chosen from some alphabet.

For example, 01101 is a string from the binary alphabet  $\Sigma = \{0, 1\}$

- **Empty String**

The empty string is the string with zero occurrence of symbol. It is denoted by  $\epsilon$

- **Length of String**

The number of positions in a string is called its length. The notation for the length of a string  $S$  is  $|S|$ . For example,  $|01101|=5$

#### 2.1.4.3 Language

A language is a specific set of strings over some fixed alphabet  $\Sigma$

**Example:**

- $\phi$  the empty set is language.
- $\{\epsilon\}$  the set containing empty string is a language
- The set of well-formed c program is a language.
- The set of all possible identifier is a language.

### Operation on languages:

- Concatnation:  $L_1 L_2 = \{s_1 s_2 | s_1 \in L_1 \text{ and } s_2 \in L_2\}$
- Union:  $L_1 \cup L_2 = \{s | s \in L_1 \text{ or } s \in L_2\}$
- Exponential:  $L^0 = \{\epsilon\}$   $L^1 = L$   $L^2 = LL$
- Kleen Closure:  $L^* = \bigcup_{i=0}^{\infty} L^i$

- Positive Closure  $L^+ = \bigcup_{i=1}^{\infty} L^i$

#### 2.1.4.4 Regular Expression

- Regular expression are the algebraic expressions that are used to describe token of a programming language.
- Let  $\Sigma$  be an alphabet, the regular expression over the alphabet  $\Sigma$  are defined inductively as follows:
  - **Basic Step:**
  - $\phi$  is a regular expression representing the empty language i.e  $L(\phi) = \phi$
  - $\epsilon$  is a regular expression representing the language of empty string  $\{\epsilon\}$  i.e;  $L(\epsilon) = \{\epsilon\}$
  - If 'a' is a symbol in  $\Sigma$ , then 'a' is a regular expression representing the language  $\{a\}$  i.e;  $L(a) = \{a\}$

#### Example:

- $1(1+0)^*0$  denotes the language of all string that begins with a '1' and ends with a '0'.
- $(1+0)^*00$  denotes the language of all strings that ends with 00.

#### 2.1.4.5 Regular Definition

To write regular expression for some languages can be difficult because their regular expressions can be quite complex. In those cases, we may use *regular definitions*. The regular definition is a sequence of definitions of the form,

$$\begin{aligned} d_1 &\rightarrow r_1 \\ d_2 &\rightarrow r_2 \\ &\dots\dots\dots \\ d_n &\rightarrow r_n \end{aligned}$$

Where  $d_i$  is a distinct name and  $r_i$  is a regular expression over symbols in  $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$

Where,  $\Sigma$  = Basic symbol and  $\{d_1, d_2, \dots, d_{i-1}\}$  = previously defined names.

#### Example:

Regular definition for specifying identifiers in a programming language like C.

letter  $\rightarrow A \mid B \mid C \mid \dots\dots\dots \mid Z \mid a \mid b \mid c \mid \dots\dots\dots \mid z$

underscore  $\rightarrow \_$

digit  $\rightarrow 0 \mid 1 \mid 2 \mid \dots\dots\dots \mid 9$

id  $\rightarrow (\text{letter} \mid \text{underscore}). (\text{letter} \mid \text{underscore} \mid \text{digit})^*$

If we are trying to write the regular expression representing identifiers without using regular definition, it will be complex.

$(A | B | C | \dots | Z | a | b | c | \dots | z | \_ | ) \cdot ((A | B | C | \dots | Z | a | b | c | \dots | z | \_ | ) (1 | 2 | \dots | 9)) ^*$

### 2.1.5 Recognition of Tokens

To recognize tokens lexical analyzer performs following steps:

- A. Lexical analyzers store the input in input buffer.
- B. The token is read from input buffer and regular expressions are built for corresponding token
- C. From these regular expressions finite automata is built. Usually NFA is built.
- D. For each state of NFA, a function is designed and each input along the transitional edges corresponds to input parameters of these functions.
- E. The set of such functions ultimately create lexical analyzer program.

### 2.1.6 Finite Automata

Finite Automata (FA) is the simplest machine to recognize patterns. Finite Automata is a model of a computational system, consisting of a set of states, a set of possible inputs, and a rule to map each state to another state, or to itself, for any of the possible inputs. Formal specification of machine is  $\{Q, \Sigma, q, F, \delta\}$  where

- Q: Finite set of states,
- $\Sigma$ : set of Input Symbols,
- q: Initial state,
- F: set of Final States
- $\delta$ : Transition Function

Finite automata come in two flavors:

- non-deterministic finite automata (NFA) have no restrictions on the labels of their edges. A symbol can label several edges out of the same state, and  $\epsilon$ , the empty string, is a possible label.
- Deterministic finite automata (DFA) have, for each state, and for each symbol of its input alphabet exactly one edge with that symbol leaving that state.



### 2.1.6.1 DFA

In DFA, for each input symbol, one can determine the state to which the machine will move. Hence, it is called Deterministic Automaton. As it has a finite number of states, the machine is called Deterministic Finite Machine or Deterministic Finite Automaton.

A DFA can be represented by a 5-tuple  $(Q, \Sigma, \delta, q_0, F)$  where.

- $Q$  is a finite set of states.
- $\Sigma$  is a finite set of symbols calling the alphabet.
- $\delta$  is the transition function where  $\delta: Q \times \Sigma \rightarrow Q$
- $q_0$  is the initial state from where any input is processed ( $q_0 \in Q$ ).
- $F$  is a set of final state/states of  $Q$  ( $F \subseteq Q$ ).

Example

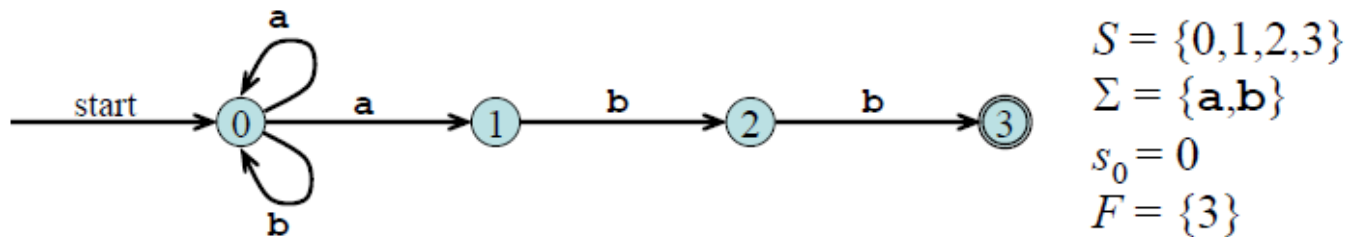


Fig: - NFA for regular expression  $(a + b)^*a b b$

### 2.1.6.2 $\epsilon$ -NFA

In NFA if a transition made without any input symbol is called  $\epsilon$ -NFA. Here we need  $\epsilon$ -NFA because the regular expressions are easily convertible to  $\epsilon$ -NFA.

A  $\epsilon$ -NFA can be represented by a 5-tuple  $(Q, \Sigma, \delta, q_0, F)$  where.

- $Q$  is a finite set of states.
- $\Sigma$  is a finite set of symbols calling the alphabet.
- $\delta$  is the transition function where  $\delta: Q \times (\Sigma \cup \epsilon) \rightarrow 2^Q$
- $q_0$  is the initial state from where any input is processed ( $q_0 \in Q$ ).
- $F$  is a set of final state/states of  $Q$  ( $F \subseteq Q$ ).

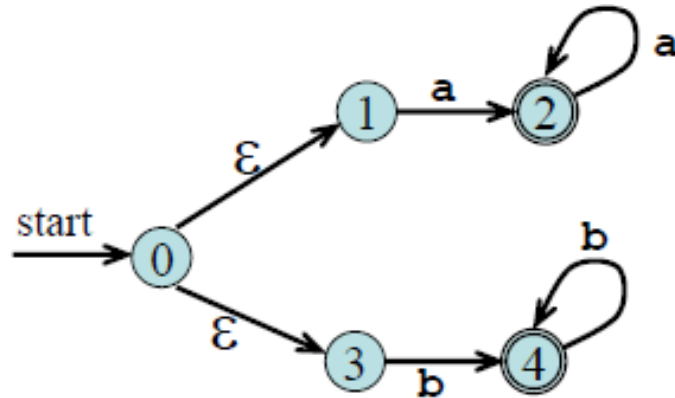


Fig: -  $\epsilon$ -NFA for regular expression  $aa^* + bb^*$

### 2.1.6.3 NFA

In NDFA, for a particular input symbol, the machine can move to any combination of the states in the machine. In other words, the exact state to which the machine moves cannot be determined. Hence, it is called Non-deterministic Automaton. As it has finite number of states, the machine is called Non-deterministic Finite Machine or Non-deterministic Finite Automaton.

A NFA can be represented by a 5-tuple  $(Q, \Sigma, \delta, q_0, F)$  where.

- $Q$  is a finite set of states.
- $\Sigma$  is a finite set of symbols calling the alphabet.
- $\delta$  is the transition function where  $\delta: Q \times \Sigma \rightarrow 2^Q$
- $q_0$  is the initial state from where any input is processed ( $q_0 \in Q$ ).
- $F$  is a set of final state/states of  $Q$  ( $F \subseteq Q$ ).

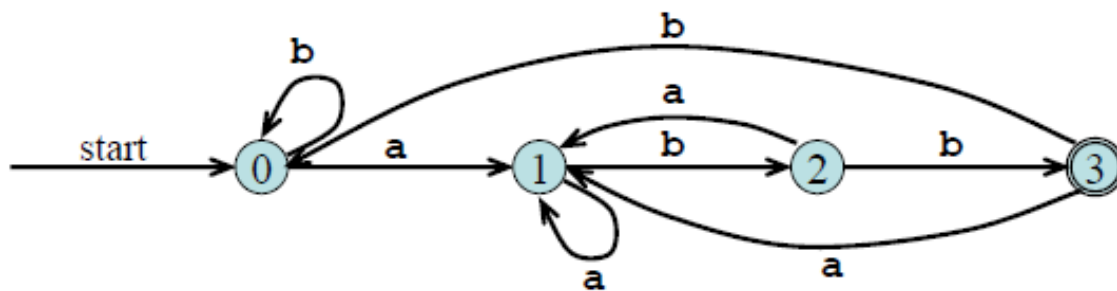
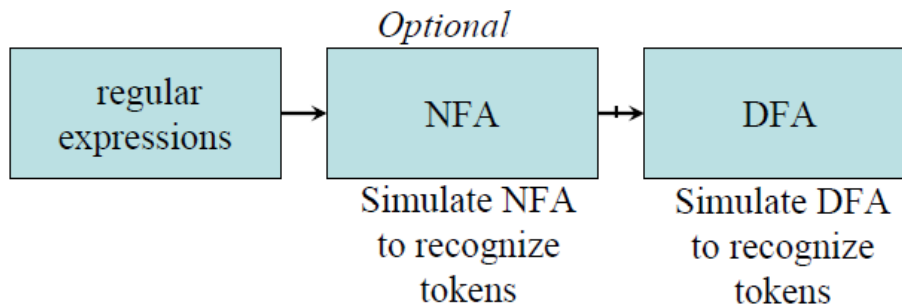


Fig:-DFA for regular expression  $(a+b)^*abb$

### 2.1.7 Design of Lexical Analyzer

First, we define regular expression for tokens, then we convert them into a DFA to get a lexical analyzer for our tokens.



Algorithm1:

- Regular Expression  $\rightarrow$  NFA  $\rightarrow$  DFA (two steps: first NFA, then to DFA) **(Self-Study)**

Algorithm2:

- Regular Expression  $\rightarrow$  DFA (Directly convert a regular expression into a DFA)

### 2.1.7.1 Conversion from RE to DFA Directly

To construct a DFA directly from a regular expression, we construct its syntax tree and then compute four functions: nullable , firstpos , lastpos , and followpos , defined as follows. Each definition refers to the syntax tree for a particular augmented regular expression ( r ) #.

- nullable ( n ) is true for a syntax-tree node n if and only if the subexpression represented by n has # in its language. That is, the subexpression can be made null or the empty string, even though there may be other strings it can represent as well.
- firstpos ( n ) is the set of positions in the subtree rooted at n that correspond to the first symbol of at least one string in the language of the subexpression rooted at n.
- lastpos ( n ) is the set of positions in the subtree rooted at n that correspond to the last symbol of at least one string in the language of the subexpression rooted at n.
- followpos ( p ), for a position p , is the set of positions q in , the entire syntax # tree such that there is some string  $x = a_1 a_2 \dots a_n$  in  $L ( r ) \#$  , such # that for some i , there is a way to explain the membership of x in  $L ( r ) \#$  by matching  $a_i$  to position p of the syntax tree and  $a_{i+1}$  to position q.

#### Conversion steps:

- Augment the given regular expression by concatenating it with special symbol # i.e.  $r \rightarrow (r) \#$
- Create the syntax tree for this augmented regular expression

In this syntax tree, all alphabet symbols (plus # and the empty string) in the augmented regular expression will be on the leaves, and all inner nodes will be the operators in that augmented regular expression.

- c) Then each alphabet symbol (plus #) will be numbered (position numbers)
- d) Traverse the tree to construct functions *nullable*, *firstpos*, *lastpos*, and *followpos*
- e) Finally construct the DFA from the *followpos*

**Rules for calculating nullable, firstpos and lastpos:**

node $n$	$\text{nullable}(n)$	$\text{firstpos}(n)$	$\text{lastpos}(n)$
is leaf labeled $\epsilon$	true	$\Phi$	$\Phi$
is leaf labeled with position $i$	false	$\{i\}$ (position of leaf node)	$\{i\}$
$n$ $\swarrow \searrow$ $c_1 \quad c_2$	$\text{nullable}(c_1)$ or $\text{nullable}(c_2)$	$\text{firstpos}(c_1) \cup \text{firstpos}(c_2)$	$\text{lastpos}(c_1) \cup \text{lastpos}(c_2)$
$n$ $\bullet$ $\swarrow \searrow$ $c_1 \quad c_2$	$\text{nullable}(c_1)$ and $\text{nullable}(c_2)$	<b>if</b> ( $\text{nullable}(c_1)$ ) <b>then</b> $\text{firstpos}(c_1) \cup \text{firstpos}(c_2)$ <b>else</b> $\text{firstpos}(c_1)$	<b>if</b> ( $\text{nullable}(c_2)$ ) <b>then</b> $\text{lastpos}(c_1) \cup \text{lastpos}(c_2)$ <b>else</b> $\text{lastpos}(c_2)$
$n$ * $\downarrow$ $c_1$	true	$\text{firstpos}(c_1)$	$\text{lastpos}(c_1)$

**Algorithm to evaluate followpos**

```

for each node  $n$  in the tree do
    if  $n$  is a cat-node with left child  $c_1$  and right child  $c_2$  then
        for each  $i$  in  $\text{lastpos}(c_1)$  do
             $\text{followpos}(i) := \text{followpos}(i) \cup \text{firstpos}(c_2)$ 
        end do
    else if  $n$  is a star-node
        for each  $i$  in  $\text{lastpos}(n)$  do
             $\text{followpos}(i) := \text{followpos}(i) \cup \text{firstpos}(n)$ 
        end do
    end if
end do

```

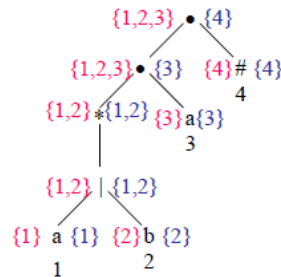
### Example

Convert regular expression  $(a \mid b)^* a$  into DFA.

Step-1: Its augmented regular expression is.

$(a|b)^*a\#$

Step-2: Construct syntax tree and calculate firstpos and followpos:



**Step-3:** Now we calculate followpos,

$\text{followpos}(1) = \{1,2,3\}$   
 $\text{followpos}(2) = \{1,2,3\}$   
 $\text{followpos}(3) = \{4\}$   
 $\text{followpos}(4) = \{\}$

Step-4: Compute State of a DFA with a help of followpos.

$S_1 = \text{firstpos}(\text{root}) = \{1,2,3\}$

mark  $S_1$

for a:  $\text{followpos}(1) \cup \text{followpos}(3) = \{1,2,3,4\} = S_2$        $\text{move}(S_1, a) = S_2$

for b:  $\text{followpos}(2) = \{1,2,3\} = S_1$        $\text{move}(S_1, b) = S_1$

mark  $S_2$

for a:  $\text{followpos}(1) \cup \text{followpos}(3) = \{1,2,3,4\} = S_2$        $\text{move}(S_2, a) = S_2$

for b:  $\text{followpos}(2) = \{1,2,3\} = S_1$        $\text{move}(S_2, b) = S_1$

Here Start state is  $S_1$  and the final state is  $S_2$ .

Step-5: Design DFA

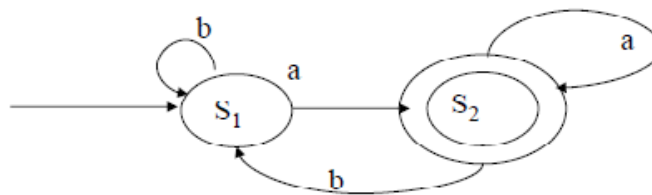


Fig: Resulting DFA of given regular expression

### 2.1.8 State Minimization in DFA

DFA minimization is the task of transforming a given deterministic finite automaton (DFA) into an equivalent DFA that has a minimum number of states. Here, two DFAs are called equivalent if they

recognize the same regular language. Several different algorithms accomplishing this task are known and described in standard textbooks on automata theory. For each regular language, there also exists a minimal automaton that accepts it, that is, a DFA with a minimum number of states and this DFA is unique (except that states can be given different names). The minimal DFA ensures minimal computational cost for tasks such as pattern matching.

Partition the set of states into two groups:

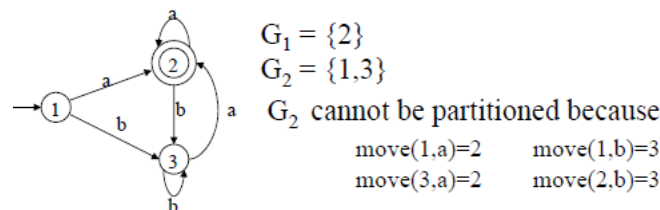
- $G_1$ : set of accepting states
- $G_2$ : set of non-accepting states

**For each new group  $G$ :**

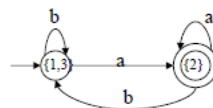
partition  $G$  into subgroups such that states  $s_1$  and  $s_2$  are in the same group if for all input symbols  $a$ , states  $s_1$  and  $s_2$  have transitions to states in the same group.

- Start state of the minimized DFA is the group containing the start state of the original DFA.
- Accepting states of the minimized DFA are the groups containing the accepting states of the original DFA.

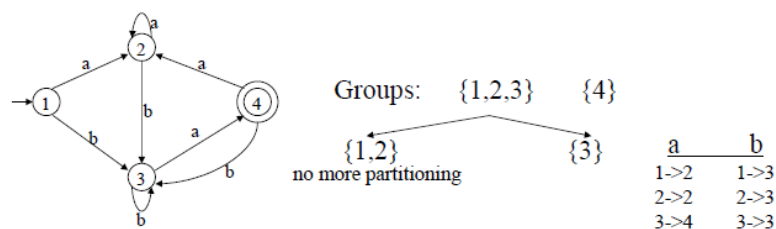
Example



So, the minimized DFA (with minimum states)



Example 2:



So, the minimized DFA (with minimum states)

