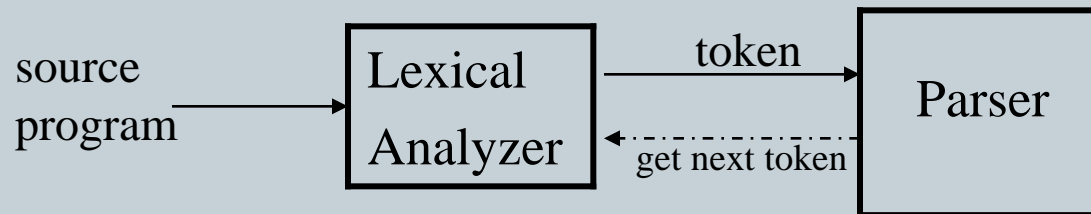


# Lexical Analyzer

1

- Lexical Analyzer reads the source program character by character to produce tokens.
- Normally a lexical analyzer doesn't return a list of tokens at one shot, it returns a token when the parser asks a token from it.



# Token

2

- Token represents a set of strings described by a pattern.
  - Identifier represents a set of strings which start with a letter continues with letters and digits
  - The actual string (newval) is called as lexeme.
  - Tokens: identifier, number, addop, delimiter, ...
- Since a token can represent more than one lexeme, additional information should be held for that specific lexeme. This additional information is called as the attribute of the token.
- For simplicity, a token may have a single attribute which holds the required information for that token.
  - For identifiers, this attribute a pointer to the symbol table, and the symbol table holds the actual attributes for that token.
- Some attributes:
  - <id,attr>                      where attr is pointer to the symbol table
  - <assgop,\_>      no attribute is needed (if there is only one assignment operator)
  - <num,val>                      where val is the actual value of the number.
- Token type and its attribute uniquely identifies a lexeme.
- Regular expressions are widely used to specify patterns.
-

# Terminology of Languages

3

- Alphabet : a finite set of symbols (ASCII characters)
- String :
  - Finite sequence of symbols on an alphabet
  - Sentence and word are also used in terms of string
  - $\epsilon$  is the empty string
  - $|s|$  is the length of string  $s$ .
- Language: sets of strings over some fixed alphabet
  - $\emptyset$  the empty set is a language.
  - $\{\epsilon\}$  the set containing empty string is a language
  - The set of well-formed C programs is a language
  - The set of all possible identifiers is a language.
- Operators on Strings:
  - Concatenation:  $xy$  represents the concatenation of strings  $x$  and  $y$ .  $s\epsilon = s$   
 $\epsilon s = s$
  - $s^n = s s s \dots s$  (  $n$  times )     $s^0 = \epsilon$

# Operations on Languages

4

- **Concatenation:**

- $L_1 L_2 = \{ s_1 s_2 \mid s_1 \in L_1 \text{ and } s_2 \in L_2 \}$

- **Union**

- $L_1 \cup L_2 = \{ s \mid s \in L_1 \text{ or } s \in L_2 \}$

- **Exponentiation:**

- $L^0 = \{\varepsilon\} \quad L^1 = L \quad L^2 = LL$

- **Kleene Closure**

- $L^* = \bigcup_{i=0}^{\infty} L^i$

- **Positive Closure**

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

# Example

5

- $L_1 = \{a,b,c,d\}$        $L_2 = \{1,2\}$
- $L_1L_2 = \{a1,a2,b1,b2,c1,c2,d1,d2\}$
- $L_1 \cup L_2 = \{a,b,c,d,1,2\}$
- $L_1^3 =$  all strings with length three (using a,b,c,d)
- $L_1^* =$  all strings using letters a,b,c,d and empty string
- $L_1^+ =$  doesn't include the empty string

# Regular Expressions

6

- We use regular expressions to describe tokens of a programming language.
- A regular expression is built up of simpler regular expressions (using defining rules)
- Each regular expression denotes a language.
- A language denoted by a regular expression is called as a **regular set**.

# Regular Expressions (Rules)

7

## Regular expressions over alphabet $\Sigma$

<u>Reg. Expr</u>	<u>Language it denotes</u>
$\varepsilon$	$\{\varepsilon\}$
$a \in \Sigma$	$\{a\}$
$(r_1) \mid (r_2)$	$L(r_1) \cup L(r_2)$
$(r_1)(r_2)$	$L(r_1) L(r_2)$
$(r)^*$	$(L(r))^*$
$(r)$	$L(r)$

- $(r)^+ = (r)(r)^*$
- $(r)? = (r) \mid \varepsilon$

# Regular Expressions (cont.)

8

- We may remove parentheses by using precedence rules.
  - \* highest
  - concatenation next
  - | lowest
- $ab^*|c$  means  $(a(b)^*)|(c)$
- Ex:
  - $\Sigma = \{0,1\}$
  - $0|1 \Rightarrow \{0,1\}$
  - $(0|1)(0|1) \Rightarrow \{00,01,10,11\}$
  - $0^* \Rightarrow \{\epsilon, 0, 00, 000, 0000, \dots\}$
  - $(0|1)^* \Rightarrow$  all strings with 0 and 1, including the empty string



# Regular Definitions

9

- 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*.
- We can give names to regular expressions, and we can use these names as symbols to define other regular expressions.
- A **regular definition** is a sequence of the definitions of the form:  
$$\begin{array}{ll} d_1 \rightarrow r_1 & \text{where } d_i \text{ is a distinct name and} \\ d_2 \rightarrow r_2 & r_i \text{ is a regular expression over symbols in} \\ \cdot & \Sigma \cup \{d_1, d_2, \dots, d_{i-1}\} \\ d_n \rightarrow r_n & \end{array}$$

basic symbols                      previously defined names

# Regular Definitions (cont.)

10

- **Ex: Identifiers in Pascal**

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

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

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

- If we try to write the regular expression representing identifiers without using regular definitions, that regular expression will be complex.

$(A \mid \dots \mid Z \mid a \mid \dots \mid z) ( (A \mid \dots \mid Z \mid a \mid \dots \mid z) \mid (0 \mid \dots \mid 9) )^*$

- **Ex: Unsigned numbers in Pascal**

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

digits  $\rightarrow \text{digit}^+$

opt-fraction  $\rightarrow ( . \text{digits} ) ?$

opt-exponent  $\rightarrow ( E (+|-)? \text{digits} ) ?$

unsigned-num  $\rightarrow \text{digits} \text{opt-fraction} \text{opt-exponent}$

## C-language identifiers

letter\_ [A-Za-z\_]

digit [0-9]

CId letter\_ ( letter\_ | digit )\*

## Unsigned integer or floating point numbers

digit [0-9]

digits digit+

number digits ( . digits )? ( E[+-]? digits )?

# Finite Automata

11

- A *recognizer* for a language is a program that takes a string  $x$ , and answers “yes” if  $x$  is a sentence of that language, and “no” otherwise.
- We call the recognizer of the tokens as a *finite automaton*.
- A finite automaton can be: *deterministic(DFA)* or *non-deterministic (NFA)*
- This means that we may use a deterministic or non-deterministic automaton as a lexical analyzer.
- Both deterministic and non-deterministic finite automaton recognize regular sets.
- Which one?
  - deterministic – faster recognizer, but it may take more space
  - non-deterministic – slower, but it may take less space
  - Deterministic automata are widely used lexical analyzers.
- First, we define regular expressions 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 to NFA, then to DFA)
  - Algorithm2: Regular Expression  $\rightarrow$  DFA (directly convert a regular expression into a DFA)

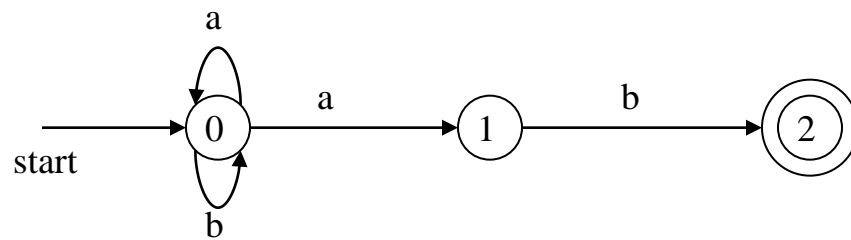
# Non-Deterministic Finite Automaton (NFA)

12

- A non-deterministic finite automaton (NFA) is a mathematical model that consists of:
  - $S$  - a set of states
  - $\Sigma$  - a set of input symbols (alphabet)
  - move – a transition function move to map state-symbol pairs to sets of states.
  - $s_0$  - a start (initial) state
  - $F$  – a set of accepting states (final states)
- $\epsilon$ - transitions are allowed in NFAs. In other words, we can move from one state to another one without consuming any symbol.
- A NFA accepts a string  $x$ , if and only if there is a path from the starting state to one of accepting states such that edge labels along this path spell out  $x$ .

# NFA (Example)

13



Transition graph of the NFA

0 is the start state  $s_0$   
 $\{2\}$  is the set of final states  $F$   
 $\Sigma = \{a, b\}$   
 $S = \{0, 1, 2\}$

Transition Function:

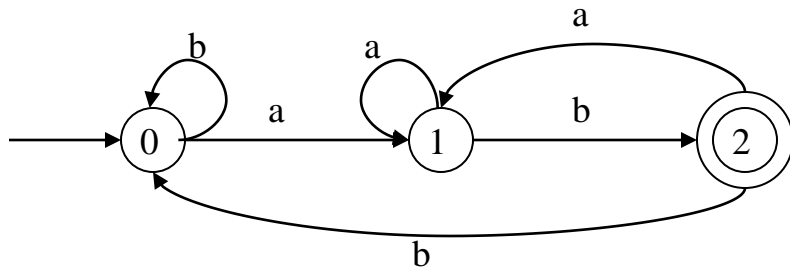
	<u>a</u>	<u>b</u>
0	$\{0, 1\}$	$\{0\}$
1	—	$\{2\}$
2	—	—

The language recognized by this NFA is  $(a|b)^* a b$

# Deterministic Finite Automaton (DFA)

14

- A Deterministic Finite Automaton (DFA) is a special form of a NFA.
  - no state has  $\epsilon$ - transition
  - for each symbol a and state s, there is at most one labeled edge a leaving s.  
i.e. transition function is from pair of state-symbol to state (not set of states)



The language recognized by  
this DFA is also  $(a|b)^* a b$

# Implementing a DFA

15

- Let us assume that the end of a string is marked with a special symbol (say eos). The algorithm for recognition will be as follows: (an efficient implementation)

```
s ← s0           { start from the initial state }
c ← nextchar       { get the next character from the input string }
while (c ≠ eos) do  { do until the end of the string }
  begin
    s ← move(s,c) { transition function }
    c ← nextchar
  end
if (s in F) then    { if s is an accepting state }
  return "yes"
else
  return "no"
```

# Implementing a NFA

16

```
S ←  $\epsilon$ -closure( $\{s_0\}$ )           { set all of states can be accessible from  $s_0$  by  $\epsilon$ -  
    transitions }  
c ← nextchar  
while (c != eos) {  
    begin  
        s ←  $\epsilon$ -closure(move(S,c)) { set of all states can be accessible from a state in S  
        c ← nextchar                by a transition on c }  
    end  
    if ( $S \cap F \neq \Phi$ ) then        { if S contains an accepting state }  
        return "yes"  
    else  
        return "no"
```

- This algorithm is not efficient.



# Converting A Regular Expression into A NFA (Thomson's Construction)

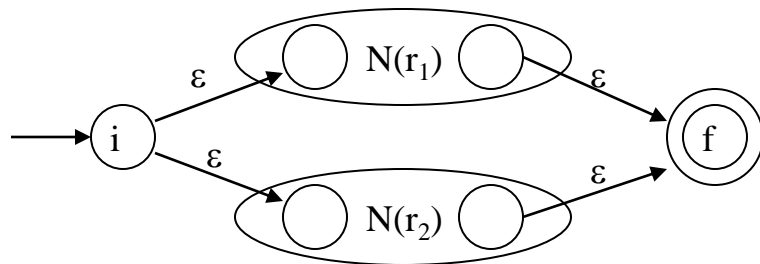
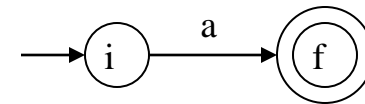
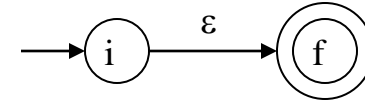
17

- This is one way to convert a regular expression into a NFA.
- There can be other ways (much efficient) for the conversion.
- Thomson's Construction is simple and systematic method. It guarantees that the resulting NFA will have exactly one final state, and one start state.
- Construction starts from simplest parts (alphabet symbols).
- To create a NFA for a complex regular expression, NFAs of its sub-expressions are combined to create its NFA,

# Thomson's Construction (cont.)

18

- To recognize an empty string  $\varepsilon$
- To recognize a symbol  $a$  in the alphabet  $\Sigma$
- If  $N(r_1)$  and  $N(r_2)$  are NFAs for regular expressions  $r_1$  and  $r_2$ 
  - For regular expression  $r_1 \mid r_2$

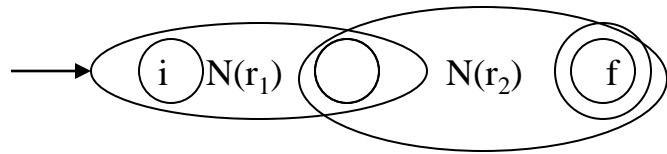


NFA for  $r_1 \mid r_2$

# Thomson's Construction (cont.)

19

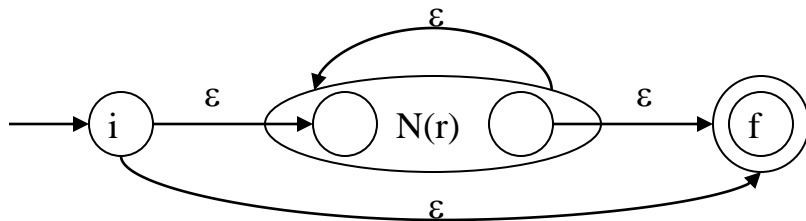
- For regular expression  $r_1 r_2$



Final state of  $N(r_2)$  become final state of  $N(r_1 r_2)$

NFA for  $r_1 r_2$

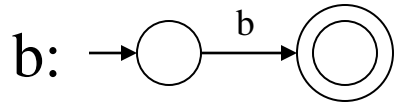
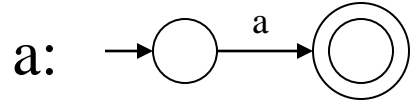
- For regular expression  $r^*$



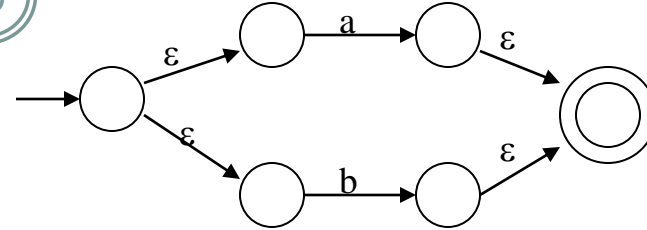
NFA for  $r^*$

# Thomson's Construction (Example - $(a|b)^* a$ )

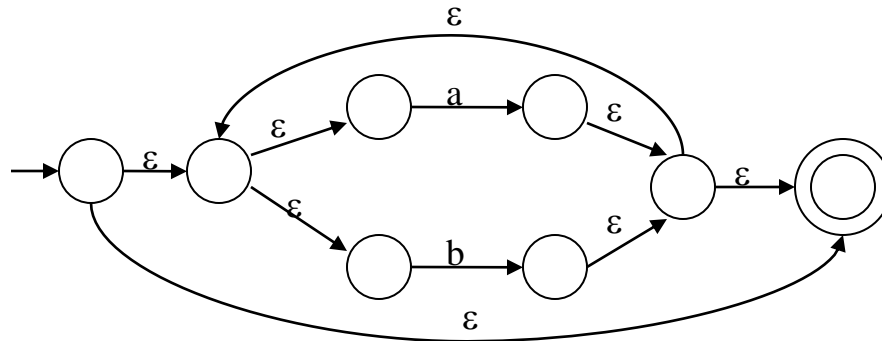
20



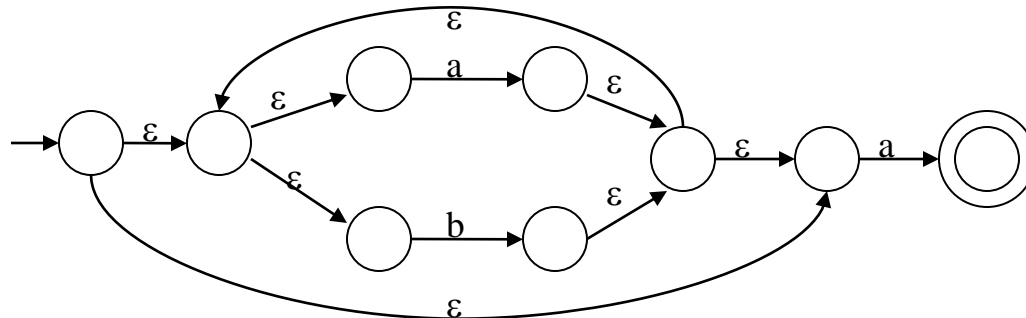
$(a | b)$



$(a|b)^*$



$(a|b)^* a$



# Converting a NFA into a DFA (subset construction)

21

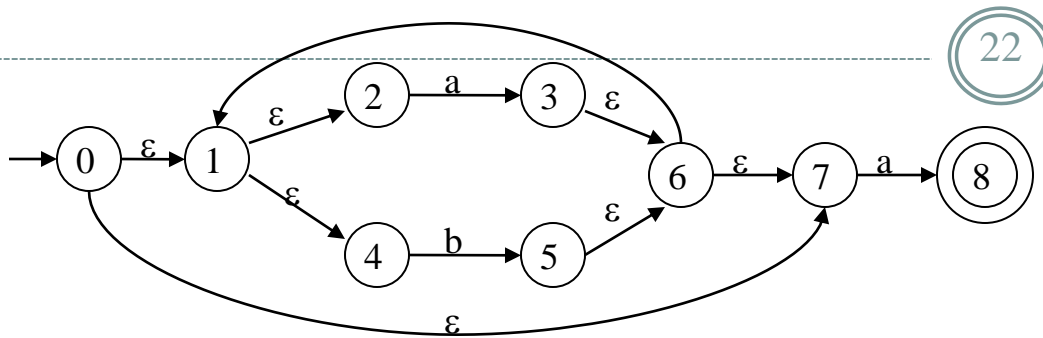
```
put  $\epsilon$ -closure( $\{s_0\}$ ) as an unmarked state into the set of DFA (DS)
while (there is one unmarked  $S_1$  in DS) do
  begin
    mark  $S_1$ 
    for each input symbol  $a$  do
      begin
         $S_2 \leftarrow \epsilon$ -closure(move( $S_1, a$ ))
        if ( $S_2$  is not in DS) then
          add  $S_2$  into DS as an unmarked state
        transfunc[ $S_1, a$ ]  $\leftarrow S_2$ 
      end
    end
  end
end
```

$\epsilon$ -closure( $\{s_0\}$ ) is the set of all states can be accessible from  $s_0$  by  $\epsilon$ -transition.

set of states to which there is a transition on  $a$  from a state  $s$  in  $S_1$

- a state  $S$  in DS is an accepting state of DFA if a state in  $S$  is an accepting state of NFA
- the start state of DFA is  $\epsilon$ -closure( $\{s_0\}$ )

# Converting a NFA into a DFA (Example)



$$S_0 = \varepsilon\text{-closure}(\{0\}) = \{0,1,2,4,7\}$$

$S_0$  into DS as an unmarked state

$\Downarrow$  mark  $S_0$

$$\varepsilon\text{-closure}(\text{move}(S_0, a)) = \varepsilon\text{-closure}(\{3,8\}) = \{1,2,3,4,6,7,8\} = S_1$$

$S_1$  into DS

$$\varepsilon\text{-closure}(\text{move}(S_0, b)) = \varepsilon\text{-closure}(\{5\}) = \{1,2,4,5,6,7\} = S_2$$

$S_2$  into DS

$$\text{transfunc}[S_0, a] \leftarrow S_1$$

$$\text{transfunc}[S_0, b] \leftarrow S_2$$

$\Downarrow$  mark  $S_1$

$$\varepsilon\text{-closure}(\text{move}(S_1, a)) = \varepsilon\text{-closure}(\{3,8\}) = \{1,2,3,4,6,7,8\} = S_1$$

$$\varepsilon\text{-closure}(\text{move}(S_1, b)) = \varepsilon\text{-closure}(\{5\}) = \{1,2,4,5,6,7\} = S_2$$

$$\text{transfunc}[S_1, a] \leftarrow S_1$$

$$\text{transfunc}[S_1, b] \leftarrow S_2$$

$\Downarrow$  mark  $S_2$

$$\varepsilon\text{-closure}(\text{move}(S_2, a)) = \varepsilon\text{-closure}(\{3,8\}) = \{1,2,3,4,6,7,8\} = S_1$$

$$\varepsilon\text{-closure}(\text{move}(S_2, b)) = \varepsilon\text{-closure}(\{5\}) = \{1,2,4,5,6,7\} = S_2$$

$$\text{transfunc}[S_2, a] \leftarrow S_1$$

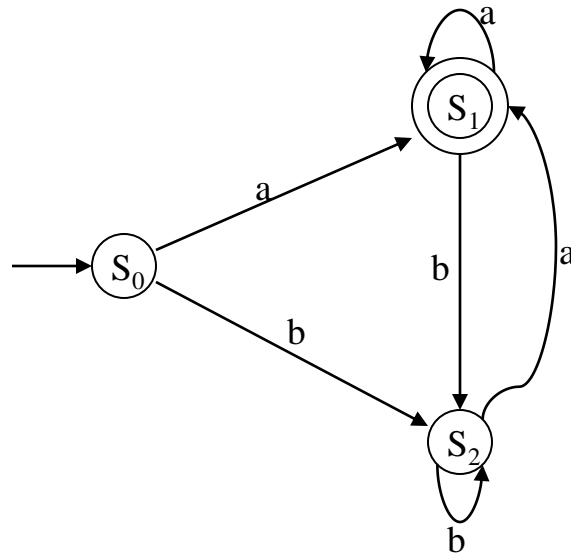
$$\text{transfunc}[S_2, b] \leftarrow S_2$$

# Converting a NFA into a DFA (Example – cont.)

23

$S_0$  is the start state of DFA since 0 is a member of  $S_0 = \{0, 1, 2, 4, 7\}$

$S_1$  is an accepting state of DFA since 8 is a member of  $S_1 = \{1, 2, 3, 4, 6, 7, 8\}$



Stop here

# Converting Regular Expressions Directly to DFAs

24

- We may convert a regular expression into a DFA (without creating a NFA first).
- First we augment the given regular expression by concatenating it with a special symbol #.  
$$r \rightarrow (r)\# \quad \text{augmented regular expression}$$
- Then, we create a 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.
- Then each alphabet symbol (plus #) will be numbered (position numbers).

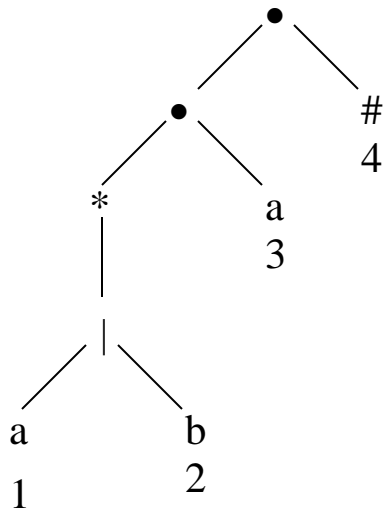


# Regular Expression $\rightarrow$ DFA (cont.)

25

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

augmented regular expression



Syntax tree of  $(a|b)^* a \#$

- each symbol is numbered (positions)
- each symbol is at a leaf
- inner nodes are operators

# followpos

26

Then we define the function **followpos** for the positions (positions assigned to leaves).

**followpos(i)** -- is the set of positions which can follow the position *i* in the strings generated by the augmented regular expression.

For example,  $(a \mid b)^* a \#$   
                  1  2      3  4

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

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

$\text{followpos}(3) = \{4\}$

$\text{followpos}(4) = \{\}$

*followpos is just defined for leaves,  
it is not defined for inner nodes.*

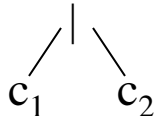
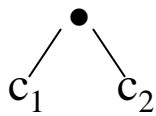

# firstpos, lastpos, nullable

27

- To evaluate followpos, we need three more functions to be defined for the nodes (not just for leaves) of the syntax tree.
- **firstpos(n)** -- the set of the positions of the **first** symbols of strings generated by the sub-expression rooted by n.
- **lastpos(n)** -- the set of the positions of the **last** symbols of strings generated by the sub-expression rooted by n.
- **nullable(n)** -- *true* if the empty string is a member of strings generated by the sub-expression rooted by n *false* otherwise

# How to evaluate firstpos, lastpos, nullable

28

<u>n</u>	<u>nullable(n)</u>	<u>firstpos(n)</u>	<u>lastpos(n)</u>
leaf labeled $\varepsilon$	true	$\Phi$	$\Phi$
leaf labeled with position i	false	$\{i\}$	$\{i\}$
	$\text{nullable}(c_1) \text{ or } \text{nullable}(c_2)$	$\text{firstpos}(c_1) \cup \text{firstpos}(c_2)$	$\text{lastpos}(c_1) \cup \text{lastpos}(c_2)$
	$\text{nullable}(c_1) \text{ and } \text{nullable}(c_2)$	if ( $\text{nullable}(c_1)$ ) $\text{firstpos}(c_1) \cup \text{firstpos}(c_2)$ else $\text{firstpos}(c_1)$	if ( $\text{nullable}(c_2)$ ) $\text{lastpos}(c_1) \cup \text{lastpos}(c_2)$ else $\text{lastpos}(c_2)$
	true	$\text{firstpos}(c_1)$	$\text{lastpos}(c_1)$

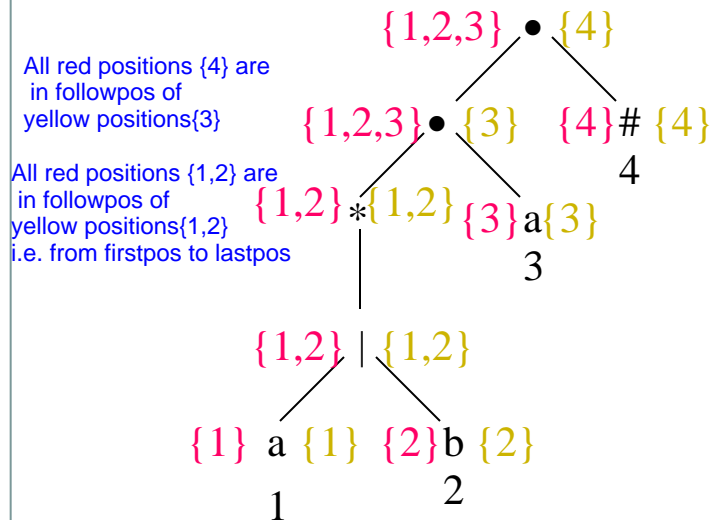
# How to evaluate followpos

29

- Two-rules define the function followpos:
  1. If **n** is concatenation-node with left child  $c_1$  and right child  $c_2$ , and **i** is a position in **lastpos**( $c_1$ ), then all positions in **firstpos**( $c_2$ ) are in **followpos**(**i**).
  2. If **n** is a star-node, and **i** is a position in **lastpos**(**n**), then all positions in **firstpos**(**n**) are in **followpos**(**i**).
- If firstpos and lastpos have been computed for each node, followpos of each position can be computed by making one depth-first traversal of the syntax tree.

# Example -- ( a | b ) \* a #

30



red— firstpos  
yellow— lastpos

Then we can calculate followpos

followpos(1) = { 1,2,3 }

followpos(2) = { 1,2,3 }

followpos(3) = { 4 }

followpos(4) = { }

- After we calculate follow positions, we are ready to create DFA for the regular expression.

# Algorithm (RE $\rightarrow$ DFA)

31

- Create the syntax tree of (r) #
- Calculate the functions: followpos, firstpos, lastpos, nullable
- Put firstpos(root) into the states of DFA as an unmarked state.
- *while* (there is an unmarked state S in the states of DFA) *do*
  - mark S
  - *for each* input symbol **a** *do*
    - ✦ let  $s_1, \dots, s_n$  are positions in S and symbols in those positions are **a**
    - ✦  $S' \leftarrow \text{followpos}(s_1) \cup \dots \cup \text{followpos}(s_n)$
    - ✦  $\text{move}(S, a) \leftarrow S'$
    - ✦ if (S' is not empty and not in the states of DFA)
      - put S' into the states of DFA as an unmarked state.
- *the start state of DFA is firstpos(root)*
- *the accepting states of DFA are all states containing the position of #*

# Example -- $(a_1 | b_2)^* a_3 \#_4$

32

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

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

$\Downarrow$  mark  $S_1$

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

$\text{move}(S_1, a) = S_2$

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

$\text{move}(S_1, b) = S_1$

$\Downarrow$  mark  $S_2$

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

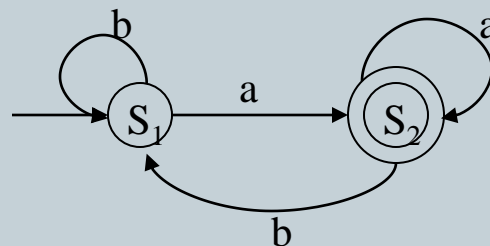
$\text{move}(S_2, a) = S_2$

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

$\text{move}(S_2, b) = S_1$

start state:  $S_1$

accepting states:  $\{S_2\}$





# Example -- ( $a_1 | \epsilon$ ) $b_2 c_3^* \#_4$

33

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

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

$\Downarrow$  mark  $S_1$

a:  $\text{followpos}(1)=\{2\}=S_2$

$\text{move}(S_1, a) = S_2$

b:  $\text{followpos}(2)=\{3,4\}=S_3$

$\text{move}(S_1, b) = S_3$

$\Downarrow$  mark  $S_2$

b:  $\text{followpos}(2)=\{3,4\}=S_3$

$\text{move}(S_2, b) = S_3$

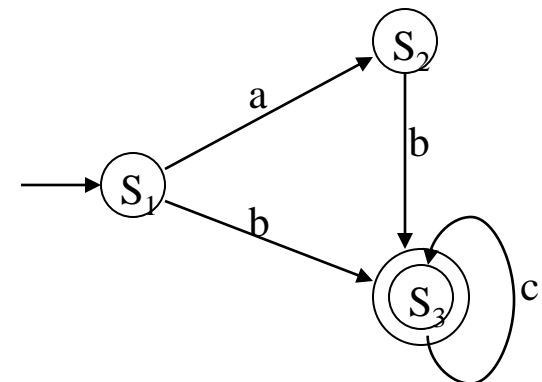
$\Downarrow$  mark  $S_3$

c:  $\text{followpos}(3)=\{3,4\}=S_3$

$\text{move}(S_3, c) = S_3$

start state:  $S_1$

accepting states:  $\{S_3\}$



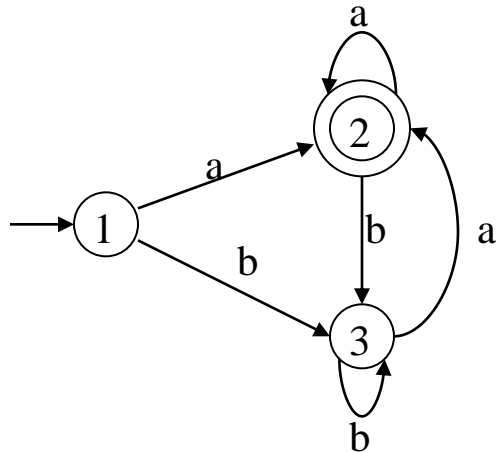
# Minimizing Number of States of a DFA

34

- 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 iff  
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.

# Minimizing DFA - Example

35



$$G_1 = \{2\}$$

$$G_2 = \{1,3\}$$

$G_2$  cannot be partitioned because

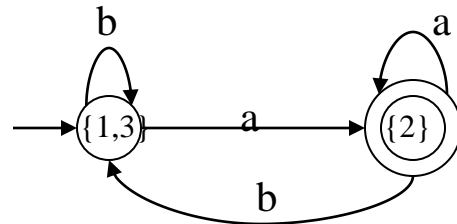
$$\text{move}(1,a)=2$$

$$\text{move}(1,b)=3$$

$$\text{move}(3,a)=2$$

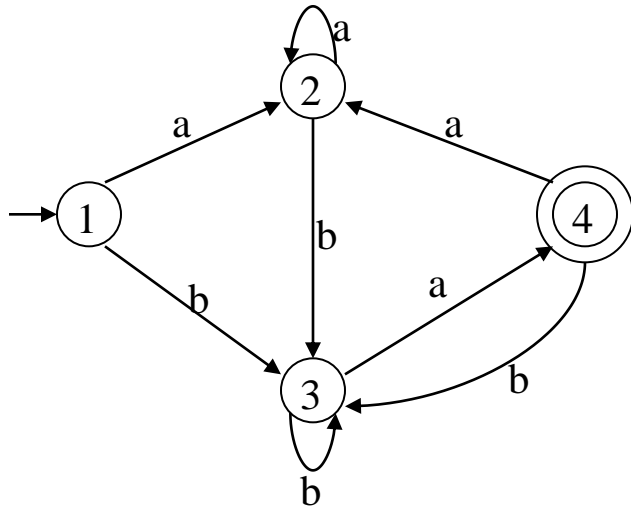
$$\text{move}(3,b)=3$$

So, the minimized DFA (with minimum states)

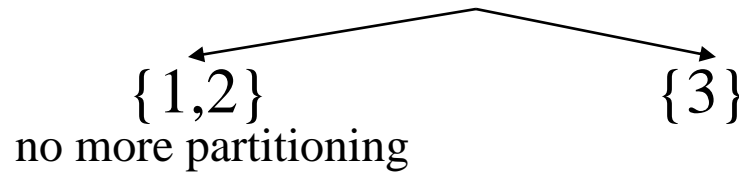


# Minimizing DFA – Another Example

36

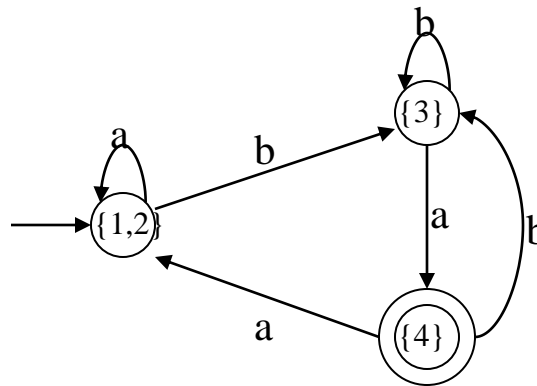


Groups: {1,2,3} {4}



a	b
1->2	1->3
2->2	2->3
3->4	3->3

So, the minimized DFA



# Some Other Issues in Lexical Analyzer

37

- The lexical analyzer has to recognize the longest possible string.
  - Ex: `identifier newval -- n ne new newv newva newval`
- What is the end of a token? Is there any character which marks the end of a token?
  - It is normally not defined.
  - If the number of characters in a token is fixed, in that case no problem: `+ -`
  - But `< ➔ <` or `<>` (in Pascal)
  - The end of an identifier : the characters cannot be in an identifier can mark the end of token.
  - We may need a lookahead
    - ✦ In Prolog: `p :- X is 1.`      `p :- X is 1.5.`  
The dot followed by a white space character can mark the end of a number.  
But if that is not the case, the dot must be treated as a part of the number.

# Some Other Issues in Lexical Analyzer (cont.)

38

- **Skipping comments**
  - Normally we don't return a comment as a token.
  - We skip a comment, and return the next token (which is not a comment) to the parser.
  - So, the comments are only processed by the lexical analyzer, and they don't complicate the syntax of the language.
- **Symbol table interface**
  - symbol table holds information about tokens (at least lexeme of identifiers)
  - how to implement the symbol table, and what kind of operations.
    - ✦ hash table – open addressing, chaining
    - ✦ putting into the hash table, finding the position of a token from its lexeme.
- **Positions of the tokens in the file (for the error handling).**