

# Compilers

by  
Marwa Yusuf

**Lecture 3**  
**Mon. 29-3-2021**

**Chapter 3 (3.3 to 3.6)**

## **Lexical Analysis**

# Specification of Tokens

- Regular expressions are used to specify lexeme patterns.

# Strings and Languages

- **Symbols:** digits, letters, punctuation.
- ***Alphabet*  $\Sigma$ :** any finite set of symbols.
  - $\{0, 1\}$  binary alphabet.
  - ASCII
  - Unicode: 100,000 symbols.
- ***String*** (over an alphabet): finite seq. of symbols drawn from alphabet.
- $|s|$ : length of string, number of occurrences of symbols.
- ***Empty string*  $\varepsilon$ :** the zero length string.

# Strings and Languages

- **Language:** any countable set of strings over some fixed alphabet, meaning is not a condition.
  - $\emptyset$ , empty set,  $\{\varepsilon\}$  are languages.
  - All syntactically well-formed c programs.
  - All grammatically correct English sentences.
- **Concatenation** of  $x$  and  $y$  ( $xy$ ): appending  $y$  to  $x$ .
  - $\varepsilon S = S\varepsilon = S$
- **Exponentiation:**
  - $S^0 = \varepsilon$
  - $S^i = S^{i-1}S$
  - $S^1 = S, S^2 = SS, S^3 = SSS \dots$

# Operations on Languages

OPERATION	DEFINITION AND NOTATION
<i>Union of <math>L</math> and <math>M</math></i>	$L \cup M = \{s \mid s \text{ is in } L \text{ or } s \text{ is in } M\}$
<i>Concatenation of <math>L</math> and <math>M</math></i>	$LM = \{st \mid s \text{ is in } L \text{ and } t \text{ is in } M\}$
<i>Kleene closure of <math>L</math></i>	$L^* = \bigcup_{i=0}^{\infty} L^i$
<i>Positive closure of <math>L</math></i>	$L^+ = \bigcup_{i=1}^{\infty} L^i$

- $L = \{A, B, \dots, Z, a, b, \dots, z\} - D = \{0, 1, \dots, 9\}$ 
  - Alphabet or a language of one letter
  - 1)  $L \cup D$
  - 2)  $LD$
  - 3)  $L^4$
  - 4)  $L^*$
  - 5)  $L(L \cup D)^*$
  - 6)  $D^+$

# Regular Expressions

- Item 5 + underscore describes C identifiers.
- ***Regular expressions***: all the languages built by applying prev. operators on some alphabet symbols.
- **Ex:** *letter\_ ( letter\_ | digit )\**
- Each regular expression  $r$  denotes a language  $L(r)$ , recursively from  $r$ 's sub-expressions.

# Regular Expressions

- **Basis:**
  - $\epsilon$  is a regular expression,  $L(\epsilon)$  is  $\{\epsilon\}$
  - If  $a$  is a symbol in  $\Sigma$  then  $\mathbf{a}$  is a regular expression and  $L(\mathbf{a}) = \{a\}$ .
- **Induction:** suppose  $r$  and  $s$  are regular expression with  $L(r)$  and  $L(s)$  languages:
  - $(r)|(s)$  is a regular expression denoting the language  $L(r) \cup L(s)$ .
  - $(r)(s)$  is a regular expression denoting the language  $L(r)L(s)$ .
  - $(r)^*$  is a regular expression denoting  $(L(r))^*$ .
  - $(r)$  is a regular expression denoting  $L(r)$ . This last rule says that we can add additional pairs of parentheses around expressions without changing the language they denote.
- Parenthesis can be dropped given that:
  - $*$ , concatenation,  $|$  with this precedence order are left associative.

# Regular Expressions

- $\Sigma = \{a, b\}$

1)  $\mathbf{a|b}$  denotes .....

2)  $\mathbf{(a|b)(a|b)}$  denotes .....

3)  $\mathbf{a^*}$  denotes .....

4)  $\mathbf{(a|b)^*}$  denotes

.....

5)  $\mathbf{a|a^*b}$  denotes .....



# Regular Expressions

- $\Sigma = \{a, b\}$ 
  - 1)  $a|b$  denotes  $\{a, b\}$
  - 2)  $(a|b)(a|b)$  denotes  $\{aa, ab, ba, bb\} = aa|ab|ba|bb$
  - 3)  $a^*$  denotes zero or more  $a = \{\epsilon, a, aa, \dots\}$
  - 4)  $(a|b)^*$  denotes all strings of zero or more  $a$  or  $b = \{\epsilon, a, b, aa, ab, bb, ba, aaa, \dots\} = (a^*b^*)^*$
  - 5)  $a|a^*b$  denotes  $\{a, b, ab, aab, aaab, \dots\} = (a)|(a^*b)$

# Regular Expressions

- ***Regular set***: a language that can be defined by a regular expression.
- If  $r$  and  $s$  denote the same regular set; they are *equivalent*.

LAW	DESCRIPTION
$r s = s r$	$ $ is commutative
$r (s t) = (r s) t$	$ $ is associative
$r(st) = (rs)t$	Concatenation is associative
$r(s t) = rs rt; (s t)r = sr tr$	Concatenation distributes over $ $
$\epsilon r = r\epsilon = r$	$\epsilon$ is the identity for concatenation
$r^* = (r \epsilon)^*$	$\epsilon$ is guaranteed in a closure
$r^{**} = r^*$	$*$ is idempotent

# Regular Definitions

- Give names to some regular expressions and use them later as if symbols.
- A regular definition:

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

...

$$d_n \rightarrow r_n$$

where:

- 1) each  $d_i$  is a new symbol not in  $\Sigma$  and not the same as any other  $d$ ;
- 2) each  $r_i$  is a regular expression over the alphabet  $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$ .

# Regular Definitions

- Ex: For C identifiers:

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

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

$id \rightarrow letter\_ (letter\_ \mid digit)^*$

- Ex: For unsigned numbers 123, 0.0334, 5.7, 1.87E-3

$digit \rightarrow$

$digits \rightarrow$

$optionalFraction \rightarrow$

$optionalExponent \rightarrow$

$number \rightarrow$

# Regular Definitions

- Ex: For C identifiers:

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

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

$id \rightarrow letter\_ ( letter\_ \mid digit )^*$

- Ex: For unsigned numbers 123, 0.0334, 5.7, 1.87E-3

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

$digits \rightarrow digit\ digit^*$

$optionalFraction \rightarrow .\ digits \mid \epsilon$

$optionalExponent \rightarrow ( E ( + \mid - \mid \epsilon ) digits ) \mid \epsilon$

$number \rightarrow digits\ optionalFraction\ optionalExponent$

# Extensions of Regular Expressions

- One or more instances  $(r)^+$ 
  - the same precedence and associativity as  $*$
  - $r^* = r^+|\epsilon$  ,  $r^+ = rr^* = r^*r$
- Zero or one instance  $r? = r|\epsilon$  [ $L(r?) = L(r) \cup L(\epsilon)$ ]
  - the same precedence and associativity as  $*$
- Character classes
  - $a_1|a_2|\dots|a_n = [a_1a_2\dots a_n]$
  - if a's form a logical sequence, like all lowercase letters, all digits, can be replaced by  $a_1$ - $a_n$
  - $[abc] = a|b|c$        $[a-z] = a|b|\dots|z$

# Extensions of Regular Expressions

- Ex: For C identifiers:

$$\textit{letter\_} \rightarrow A \mid B \mid \dots \mid Z \mid a \mid b \mid \dots \mid z \mid \_$$
$$\textit{digit} \rightarrow 0 \mid 1 \mid \dots \mid 9$$
$$\textit{id} \rightarrow \textit{letter\_} ( \textit{letter\_} \mid \textit{digit} )^*$$

- becomes

$$\textit{letter\_} \rightarrow$$
$$\textit{digit} \rightarrow$$
$$\textit{id} \rightarrow$$

# Extensions of Regular Expressions

- Ex: For C identifiers:

$$letter\_ \rightarrow A \mid B \mid \dots \mid Z \mid a \mid b \mid \dots \mid z \mid \_$$
$$digit \rightarrow 0 \mid 1 \mid \dots \mid 9$$
$$id \rightarrow letter\_ ( letter\_ \mid digit )^*$$

- becomes

$$letter\_ \rightarrow [A-Za-z]$$
$$digit \rightarrow [0-9]$$
$$id \rightarrow letter\_ ( letter\_ \mid digit )^*$$



# Regular Definitions

- Ex: For unsigned numbers 123, 0.0334, 5.7, 1.87E-3

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

*digits*  $\rightarrow \textit{digit digit}^*$

*optionalFraction*  $\rightarrow . \textit{digit} \mid \varepsilon$

*optionalExponent*  $\rightarrow ( E ( + \mid - \mid \varepsilon ) \textit{digits} ) \mid \varepsilon$

*number*  $\rightarrow \textit{digits optionalFraction optionalExponent}$

- becomes

*digit*  $\rightarrow$

*digits*  $\rightarrow$

*optionalFraction*  $\rightarrow$

*optionalExponent*  $\rightarrow$

*number*  $\rightarrow$

# Regular Definitions

- Ex: For unsigned numbers 123, 0.0334, 5.7, 1.87E-3

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

*digits*  $\rightarrow digit\ digit^*$

*optionalFraction*  $\rightarrow .\ digit \mid \epsilon$

*optionalExponent*  $\rightarrow ( E ( + \mid - \mid \epsilon )\ digits ) \mid \epsilon$

*number*  $\rightarrow digits\ optionalFraction\ optionalExponent$

- becomes

*digit*  $\rightarrow [0-9]$

*digits*  $\rightarrow digit^+$

~~*optionalFraction*  $\rightarrow .\ digit \mid \epsilon$~~

~~*optionalExponent*  $\rightarrow ( E ( + \mid - \mid \epsilon )\ digits ) \mid \epsilon$~~

*number*  $\rightarrow digits\ (. digits)?\ ( E [+ -]? digits )?$

# Recognition of Tokens

- Build a pattern matching code.
- **Ex:** Grammar of **if statement** like Pascal (**then** is explicit, = and <> are for comparison).

$$\begin{aligned} stmt \rightarrow & \text{if } expr \text{ then } stmt \\ & | \text{if } expr \text{ then } stmt \text{ else } stmt \\ & | \varepsilon \end{aligned}$$
$$\begin{aligned} expr \rightarrow & term \text{ relop } term \\ & | term \end{aligned}$$
$$\begin{aligned} term \rightarrow & id \\ & | number \end{aligned}$$

# Recognition of Tokens

- The terminals are: **if**, **then**, **else**, **relop**, **number** and **id**

*if*  $\rightarrow$  if

*then*  $\rightarrow$  then

*else*  $\rightarrow$  else

*relop*  $\rightarrow$  < | > | <= | >= | = | <>

*digit*  $\rightarrow$  [0-9]

*digits*  $\rightarrow$  *digit*<sup>+</sup>

*number*  $\rightarrow$  *digits* ( . *digits*)? ( E [+ -]? *digits*)?

*letter*  $\rightarrow$  [A-Za-z]

*id*  $\rightarrow$  *letter* ( *letter* | *digit* )\*

# Recognition of Tokens

- For stripping whitespace:

$$ws \rightarrow ( \text{blank} \mid \text{tab} \mid \text{newline} )^+$$

- **blank**, **tab** and **newline** are symbols for the ASCII characters.
- *ws* is a token that is NOT returned to the parser.

# Recognition of Tokens

LEXEMES	TOKEN NAME	ATTRIBUTE VALUE
Any <i>ws</i>	-	-
<b>if</b>	<b>if</b>	-
<b>then</b>	<b>then</b>	-
<b>else</b>	<b>else</b>	-
Any <i>id</i>	<b>id</b>	Pointer to table entry
Any <i>number</i>	<b>number</b>	Pointer to table entry
<	<b>relop</b>	LT
<=	<b>relop</b>	LE
=	<b>relop</b>	EQ
<>	<b>relop</b>	NE
>	<b>relop</b>	GT
>=	<b>relop</b>	

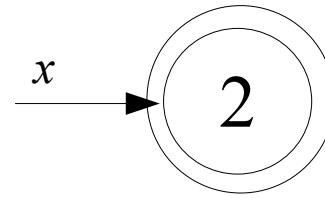
# Transition Diagrams

- Convert patterns into transition diagrams (Intermediate step).
- For now, manually.
- ***States***: (circles or nodes) a condition during scanning.
- ***Edges***: directed between states, labeled by symbol(s), *forward* pointer advances according to input and edges.
- For now, assume all transition diagrams are deterministic.

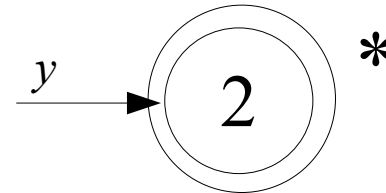
# Transition Diagrams

- Some conventions:

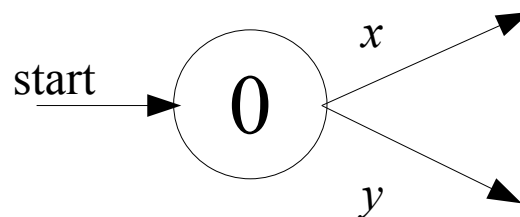
1) Certain states are **accepting** or **final** (double circle): a lexeme has been found, attached to it an action (usually return token to parser).



2) If lexeme does not include the symbol that got us to the accepting state, put a \* (or more) near the accepting state.



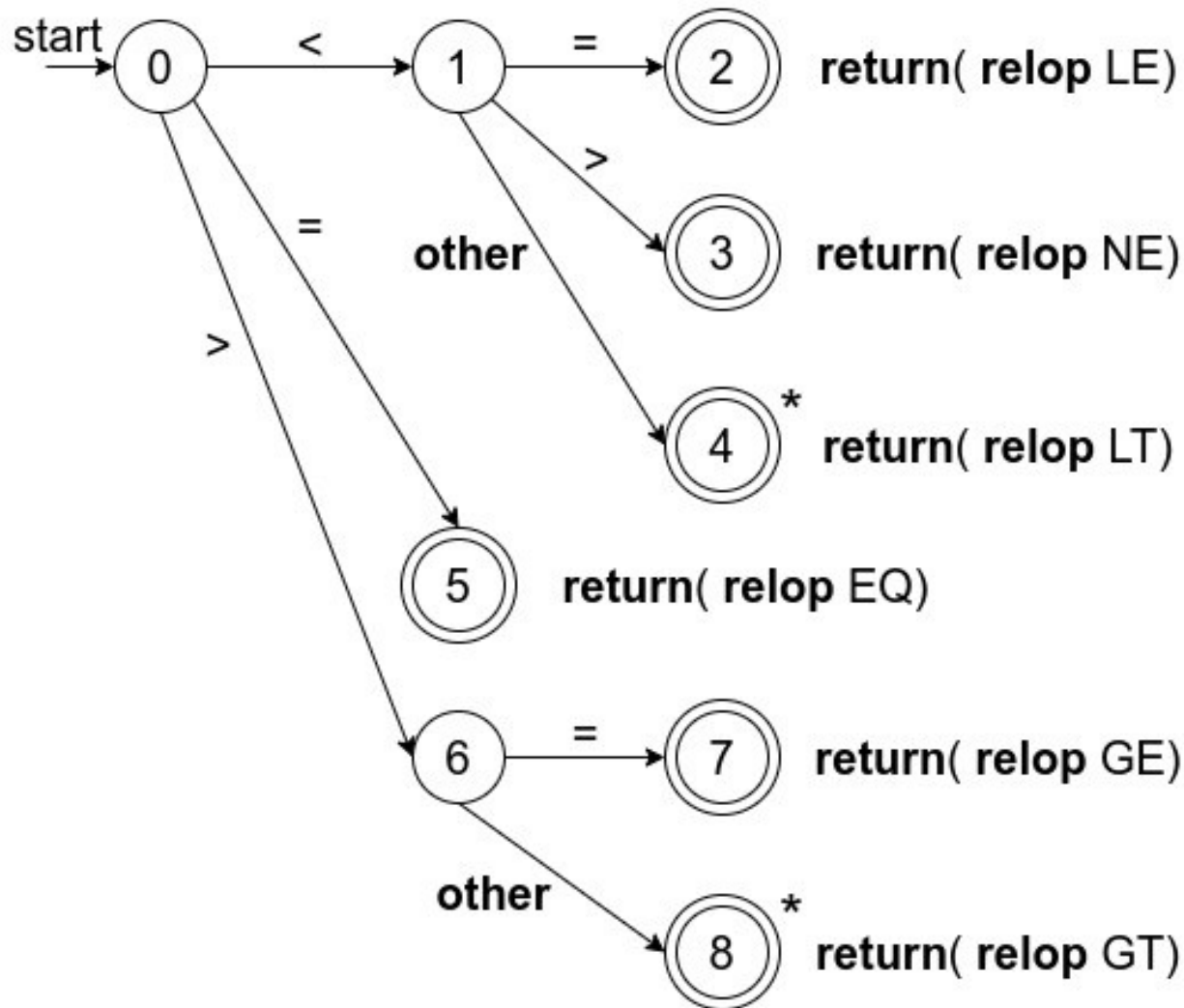
3) One state is the **start state** or **initial state**: edge labeled start coming from nowhere.





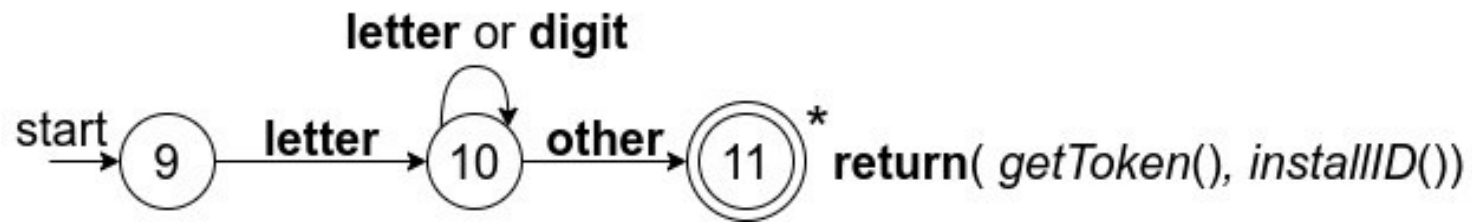
# Transition Diagrams - **relop**

*relop*  $\rightarrow$  < | > | <= | >= | = | <>



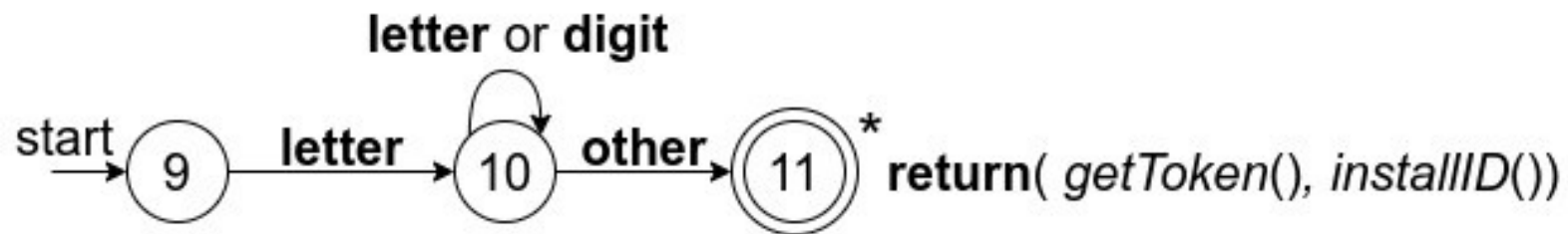
# Transition Diagrams - id

$id \rightarrow letter ( letter \mid digit )^*$

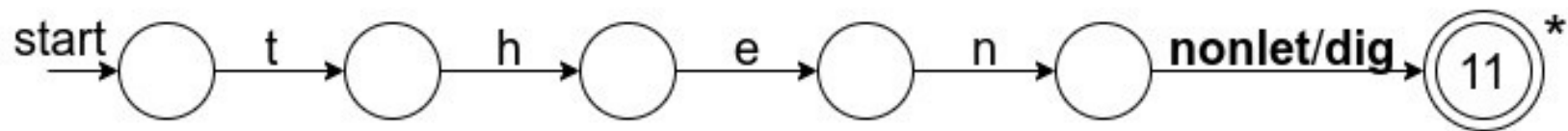


# Recognition of Reserved Words and Identifiers

- 2 Approaches to differentiate reserved words from identifiers:
  - 1) Using symbol table initially loaded with reserved words:



- 2) Create a diagram for each keyword, with higher priority for keywords. (must check for end of word, **thenextvalue**):

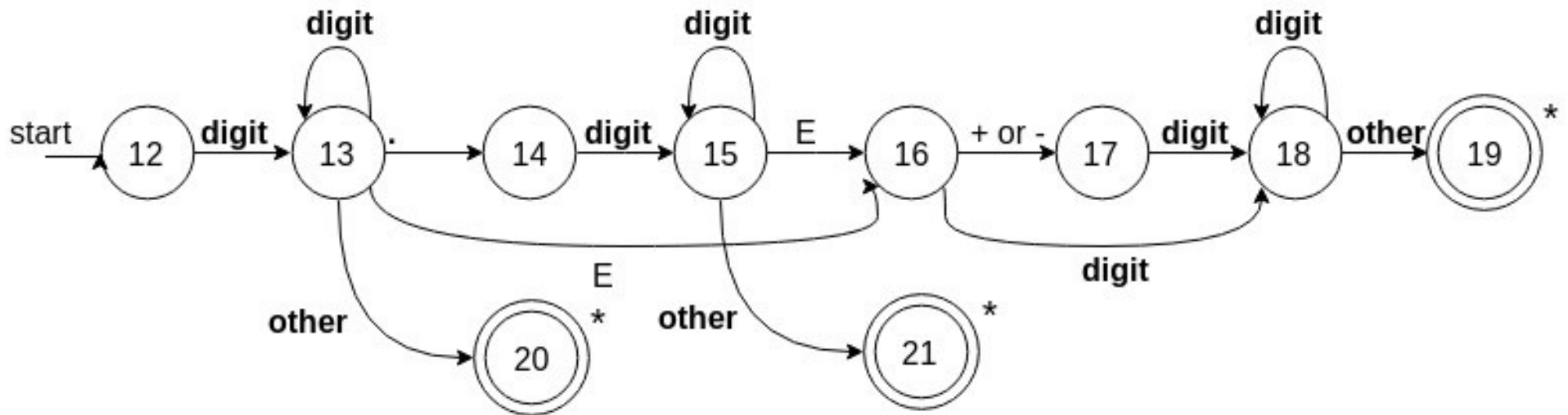


# Recognizing Unsigned Number

*digit* → [0-9]

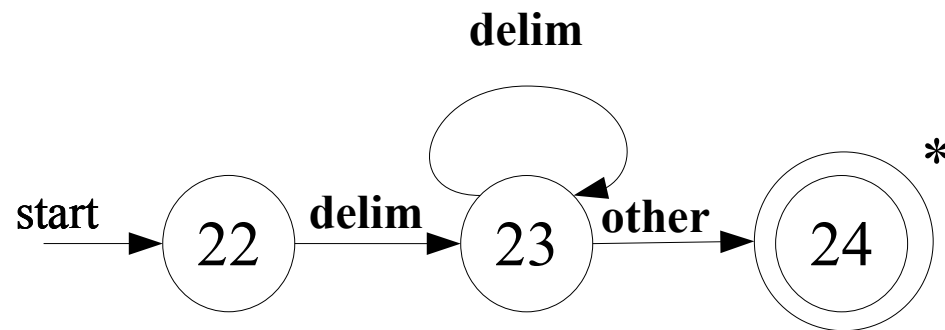
*digits* → *digit*<sup>+</sup>

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



# Recognizing Whitespace

$ws \rightarrow ( \text{blank} \mid \text{tab} \mid \text{newline} )^+$



# Transition Diagram Based Lexical Analyzer

- Each state  $\rightarrow$  a piece of code.
- A switch statement to the next state given input symbol.

# Code for **relop** Transition Diagram

```
TOKEN getRelop() {  
    TOKEN retToken = new(RELOP);  
    while(1) { /* repeat character processing until return or  
        failure */  
        switch(state) {  
            case 0: c = nextChar();  
                if ( c == '<' ) state = 1;  
                else if ( c == '=' ) state = 5;  
                else if ( c == '>' ) state = 6;  
                else fail(); /* lexeme is not a relop */  
                break;  
            case 1: ...  
            ...  
            case 8: retract();  
                retToken.attribute = GT;  
                return(retToken);  
        }  
    }  
}
```

# Transition Diagram Based Lexical Analyzer - Ways

- 1) Try diagrams sequentially, fail() resets *forward* and starts another diagram. A diagram for each keyword can be used this way, just try them before **id**.
- 2) Try diagrams in parallel. Prefer the longest prefix of the input to resolve similar prefixes.
- 3) Combine all diagrams into one. In prev. examples, combine all start states into one state. This is easy, cause they differ in first character. It is not always that easy.



# The Lexical Analyzer Generator Lex

- Skip section 3.5 (We Use ANTLR).

# Finite Automata

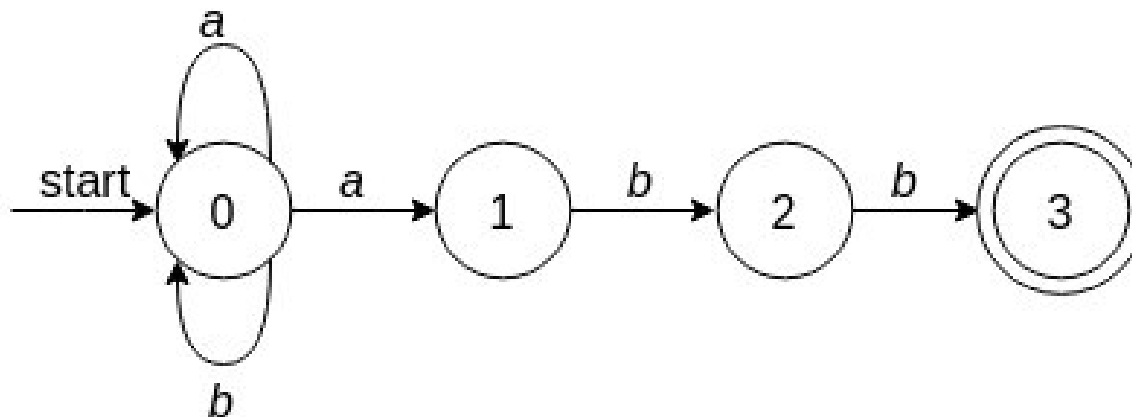
- Graphs like transition diagrams, but:
  - 1) say “yes” or “no” about an input string.
  - 2) NFA: the same symbol may label multiple edges.  
DFA: for each state and for each symbol exactly one edge.

# NFA

- Consists of:
  - 1) A finite set of states  $S$ .
  - 2) A set of input symbols  $\Sigma$ , the *input alphabet*. We assume that  $\varepsilon$ , which stands for the empty string, is never a member of  $\Sigma$ .
  - 3) A *transition function* that gives, for each state, and for each symbol in  $\Sigma \cup \{\varepsilon\}$  a set of *next states*.
  - 4) A state  $s_0$  from  $S$  that is distinguished as the *start state* (or *initial state*).
  - 5) A set of states  $F$ , a subset of  $S$ , that is distinguished as the *accepting states* (or *final states*).

# Transition Graph

- NFA can be represented by a transition graph, similar to transition diagram, except:
  - The same symbol can label multiple edges.
  - $\epsilon$  can label an edge.
- Ex:  **$(a|b)^*abb$**

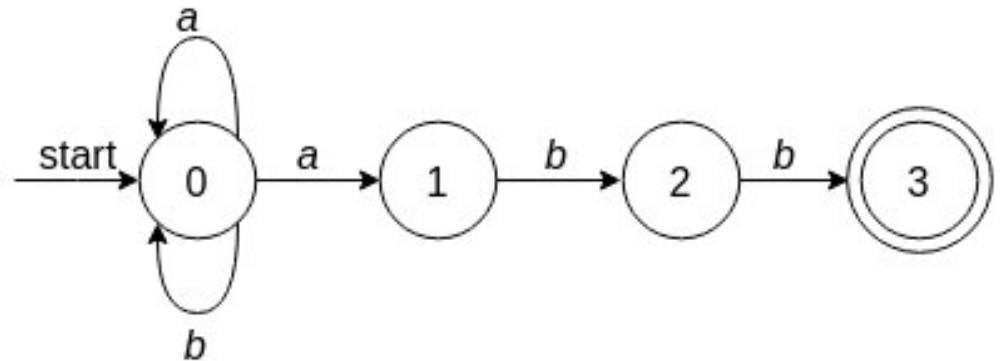


# Transition Tables

- Rows for states and columns for symbols.
- Advantage: Easy to find transition.

Disadvantage: large space with large input alphabet while most states have no move with all symbols.

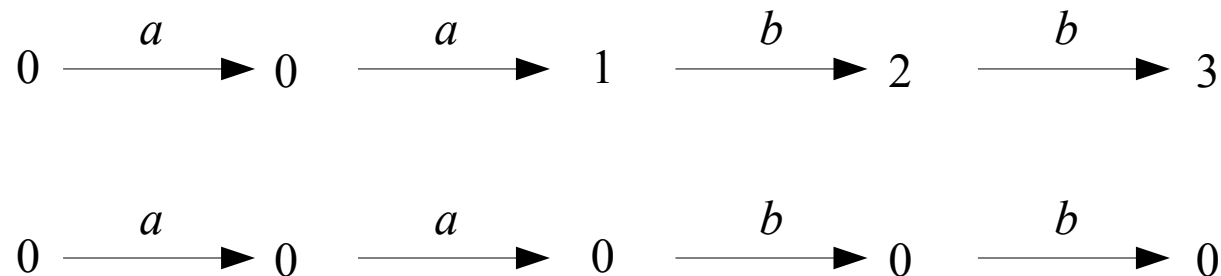
- Ex:  $(a|b)^*abb$



State	<i>a</i>	<i>b</i>	$\epsilon$
0	{0,1}	{0}	$\Phi$
1	$\Phi$	{2}	$\Phi$
2	$\Phi$	{3}	$\Phi$
3	$\Phi$	$\Phi$	$\Phi$

# Acceptance of Input Strings by Automata

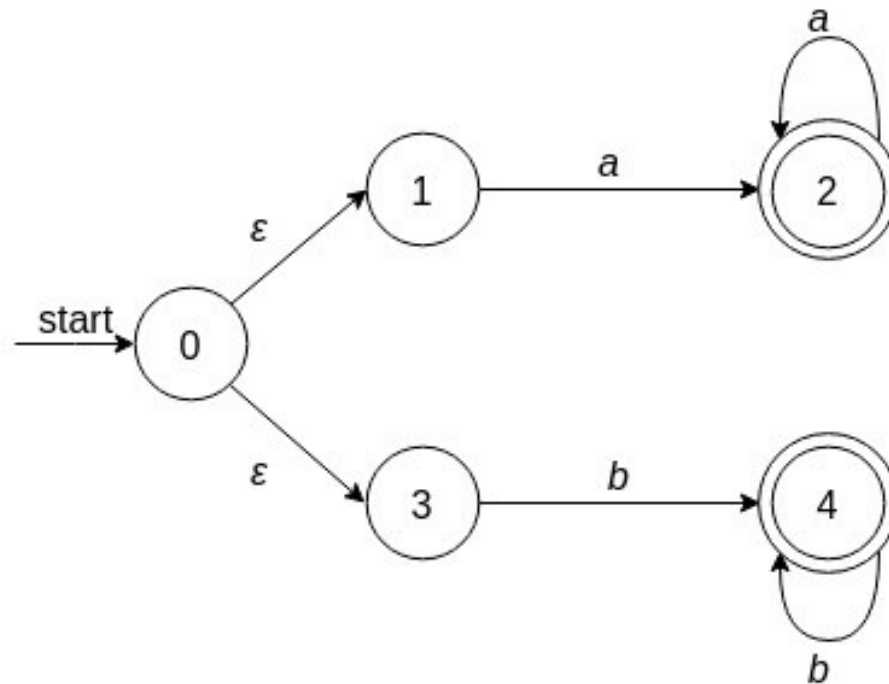
- An NFA *accepts*  $x$  iff there is some path from *start state* to one *accepting state* such that the path forms  $x$ .
- Ex:  $aabb$



- The *language defined* (or *accepted*) by an NFA is the set of strings labeling some path from the start state to an accepting state  $L(A)$ .

# Acceptance of Input Strings by Automata

- Ex:  $L(\mathbf{aa^*|bb^*})$



# DFA

- Special case of NFA where:
  - 1) There are no moves on input  $\epsilon$ , and;
  - 2) For each state  $s$  and input symbol  $a$ , there is exactly one edge out of  $s$  labeled  $a$ .
- NFA is abstract, DFA is concrete algorithm.
- Every regular expression and every NFA can be converted to DFA.
- DFA is what is implemented or simulated to build lexical analyzer.



# Simulating a DFA Algorithm

- **INPUT** : An input string  $x$  terminated by an end-of-file character **eof**. A DFA  $D$  with start state  $s_o$ , accepting states  $F$ , and transition function  $move$ .
- **OUTPUT** : Answer "yes" if  $D$  accepts  $x$ ; "no" otherwise.
- **METHOD** : Apply the following algorithm to the input string  $x$ . The function  $move(s, c)$  gives the state to which there is an edge from state  $s$  on input  $c$ . The function  $nextChar$  returns the next character of the input string  $x$ .

# Algorithm Code

$s = s_0$

$c = \text{nextChar}();$

**while** (  $c \neq \text{eof}$  ) {

$s = \text{move}(s, c);$

$c = \text{nextChar}();$

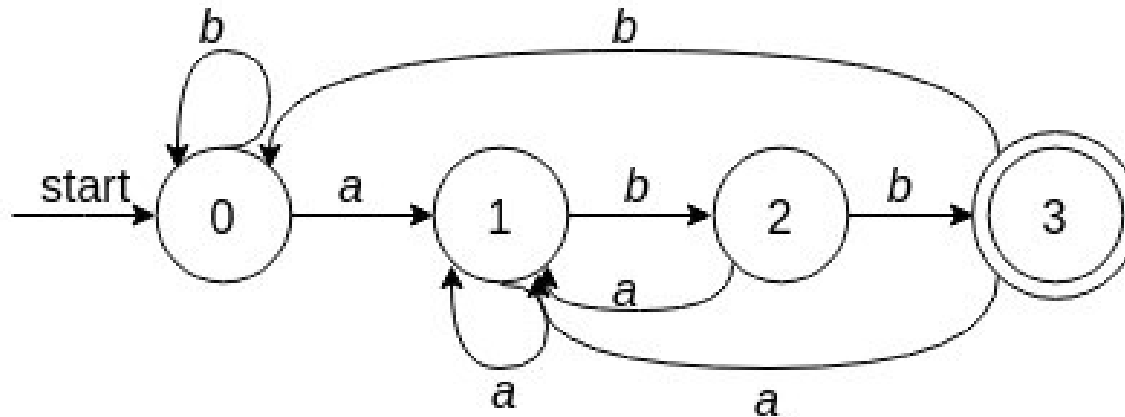
}

**if** (  $s$  is in  $F$  ) **return** "yes";

**else return** "no";

# DFA Graph

- Ex:  $(a|b)^*abb$



- Given  $ababb$ , this DFA enters the seq. 0, 1, 2, 1, 2, 3 and returns “yes”.

