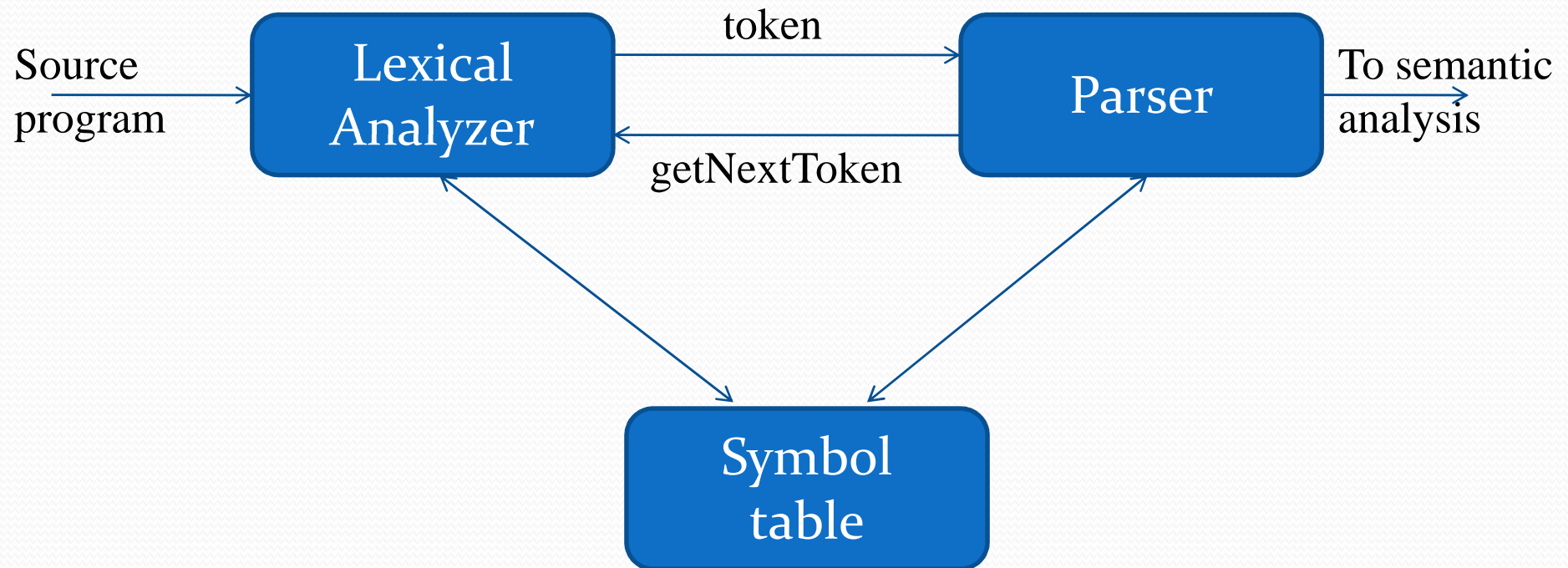




# Lexical Analyzer

## Introduction

# The Role of Lexical Analyzer



# Why to separate Lexical analysis and Parsing

## 1. Simplicity of design

Removing white space by Lexical Analyzer pays way for easy implementation of parsing

## 2. Improving compiler efficiency

Large amount of time is spend in reading the source program. Buffering techniques for reading input characters.

## 3. Enhancing compiler portability

The representation of special symbols can be isolated in Lexical Analyzer

# Tokens, Patterns and Lexemes

- A **token** is a pair of a token name and an optional token value
- A **pattern** is a description of the form that the lexemes of a token may take
- A **lexeme** is a sequence of characters in the source program that matches the pattern for a token

# Example

Token	Informal description	Sample lexemes
<b>if</b>	Characters i, f	if
<b>else</b>	Characters e, l, s, e	else
<b>comparison</b>	< or > or <= or >= or == or !=	<=, !=
<b>id</b>	Letter followed by letter and digits	pi, score, D2
<b>number</b>	Any numeric constant	3.14159, 0, 6.02e23
<b>literal</b>	Anything but “ sorrounded by “	“core dumped”

# Attributes for tokens

- $E = M * C ** 2$ 
  - $\langle \text{id, pointer to symbol table entry for } E \rangle$
  - $\langle \text{assign-op} \rangle$
  - $\langle \text{id, pointer to symbol table entry for } M \rangle$
  - $\langle \text{mult-op} \rangle$
  - $\langle \text{id, pointer to symbol table entry for } C \rangle$
  - $\langle \text{exp-op} \rangle$
  - $\langle \text{number, integer value } 2 \rangle$

# Lexical errors

- Some errors are out of power of lexical analyzer to recognize:
  - $f_i(a == f(x)) \dots$
- However it may be able to recognize errors like:
  - $d = 2r$
- Such errors are recognized when no pattern for tokens matches a character sequence

# Error recovery

- Panic mode: successive characters are ignored until we reach to a well formed token
- Delete one character from the remaining input
- Insert a missing character into the remaining input
- Replace a character by another character
- Transpose two adjacent characters





# Specification of Tokens

# Specification of tokens

- Alphabet or Character Class
  - $\Sigma$  is a finite set of symbols (characters)
  - $\{0,1\}$  is a binary alphabet
- String or Sentence or word
  - A *string*  $s$  is a finite sequence of symbols from  $\Sigma$ 
    - $|s|$  denotes the length of string  $s$
    - $\epsilon$  denotes the empty string, thus  $|\epsilon| = 0$
    - banana  $\rightarrow$   $|\text{banana}|=6$
- Language
  - A *language* is a specific set of strings over some fixed alphabet  $\Sigma$
  - $\Sigma=\{0,1\}$
  - $L=\{0,1,00,11,01,10,000,001,010,011,\dots\}$

# Specification of tokens Cont...

- **Prefix of s**
  - A string obtained by removing 0 or more trailing symbols of s
  - b, ba, ban, bana, banan, banana
- **Suffix of s**
  - A string formed by deleting 0 or more leading symbols of s
  - a, na, ana, nana ...
- **Substring of s**
  - A string obtained by removing the suffix and prefix from s
  - ana, nan etc
- **Proper prefix and Proper Suffix**
  - Any prefix or suffix other than the string itself
  - b, ba, a, nana ...
- **Subsequence of s**
  - Any string formed by deleting zero or more not necessarily contiguous symbols from s.
  - baaa, ann...

# Language Operations

- *Union*

$$L \cup M = \{s \mid s \in L \text{ or } s \in M\}$$

- *Concatenation*

$$LM = \{xy \mid x \in L \text{ and } y \in M\}$$

- *Kleene closure*

$$L^* = \cup_{i=0, \dots, \infty} L^i$$

- *Positive closure*

$$L^+ = \cup_{i=1, \dots, \infty} L^i$$

# Regular Expressions

## Rules for Regular Expression

- $\epsilon$  is a regular expression,  $L(\epsilon) = \{\epsilon\}$
- If  $a$  is a symbol in  $\Sigma$  then  $a$  is a regular expression,  $L(a) = \{a\}$
- $(r) \mid (s)$  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  $L(r)$

Ex : Identifier  $\rightarrow$  letter ( letter  $\mid$  digit ) \*

Language for Regular Expression  $\rightarrow$  Regular Set

# Precedence

- \* (Closure) has the higher precedence
- . (Concatenation) has the next higher precedence
- | (Union) has the least precedence

Remove unnecessary parentheses

$$(a) \mid ((b)^*c) \rightarrow a \mid b^*c$$

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

$$\text{RE } a \mid b \quad \{a, b\}$$

$$(a/b)(a/b) \quad \{aa, ab, ba, bb\}$$

$$(a/b)^* \quad ? \qquad a/a^*b \qquad ?$$

If 2 r.e  $r$  and  $s$  denote the same language then  $r$  and  $s$  are said to be **equivalent** ie.  $r=s$  ex.  $a/b = b/a$

# Regular definitions

- Regular definitions introduce a naming convention:

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

...

$$d_n \rightarrow r_n$$

where each  $r_i$  is a regular expression over

$$\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$$

- Any  $d_j$  in  $r_i$  can be textually substituted in  $r_i$  to obtain an equivalent set of definitions

# Regular definitions Cont...

- Example:

**letter**  $\rightarrow$  **A** | **B** | ... | **Z** | **a** | **b** | ... | **z**

**digit**  $\rightarrow$  **0** | **1** | ... | **9**

**id**  $\rightarrow$  **letter** ( **letter** | **digit** )<sup>\*</sup>

- Regular definitions are not recursive:

**digits**  $\rightarrow$  **digit digits** | **digit**      *wrong!*



# Notational Shorthand

- One or more instances:  $(r)^+$
- Zero of one instances:  $r^?$
- Character classes:  $[abc]$

$$r^+ = rr^*$$

$$r^? = r \mid \epsilon$$

$$[\mathbf{a-z}] = \mathbf{a} \mid \mathbf{b} \mid \mathbf{c} \mid \dots \mid \mathbf{z}$$

# Notational Shorthand

- letter\_  $\rightarrow$  [A-Za-z\_]
  - digit  $\rightarrow$  [0-9]
  - id  $\rightarrow$  letter\_(letter | digit)\*
- 
- Examples:  
**digit**  $\rightarrow$  [0-9]  
**num**  $\rightarrow$  **digit**<sup>+</sup> ( . **digit**<sup>+</sup> )? ( **E** ( + | - )? **digit**<sup>+</sup> )?

# Recognition of tokens

- Starting point is the language grammar to understand the tokens:

stmt  $\rightarrow$  **if** expr **then** stmt  
          | **if** expr **then** stmt **else** stmt  
          |  $\epsilon$

expr  $\rightarrow$  term **relop** term  
          | term

term  $\rightarrow$  **id**  
          | **number**

# Recognition of tokens (cont.)

- The next step is to formalize the patterns:

*digit* -> [0-9]

*digits* -> digit+

*number* -> digit(.digits)? (E[+-]? Digit)?

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

*id* -> letter (letter | digit)\*

*if* -> if

*then* -> then

*else* -> else

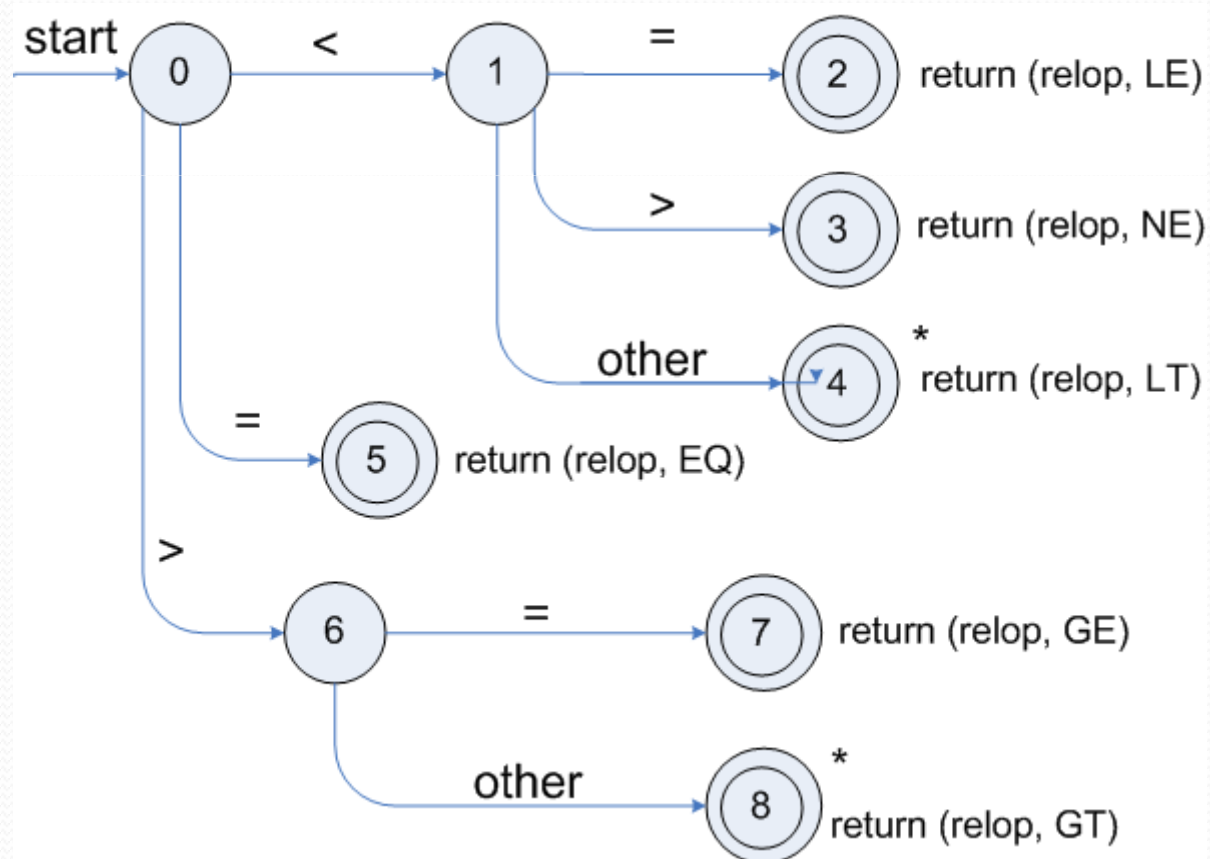
*Relop* -> < | > | <= | >= | = | <>

- We also need to handle whitespaces:

*ws* -> (blank | tab | newline)+

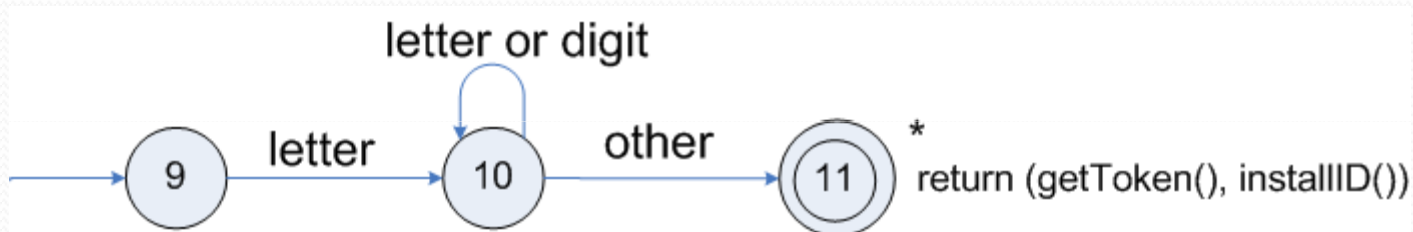
# Transition diagrams

- Transition diagram for relop



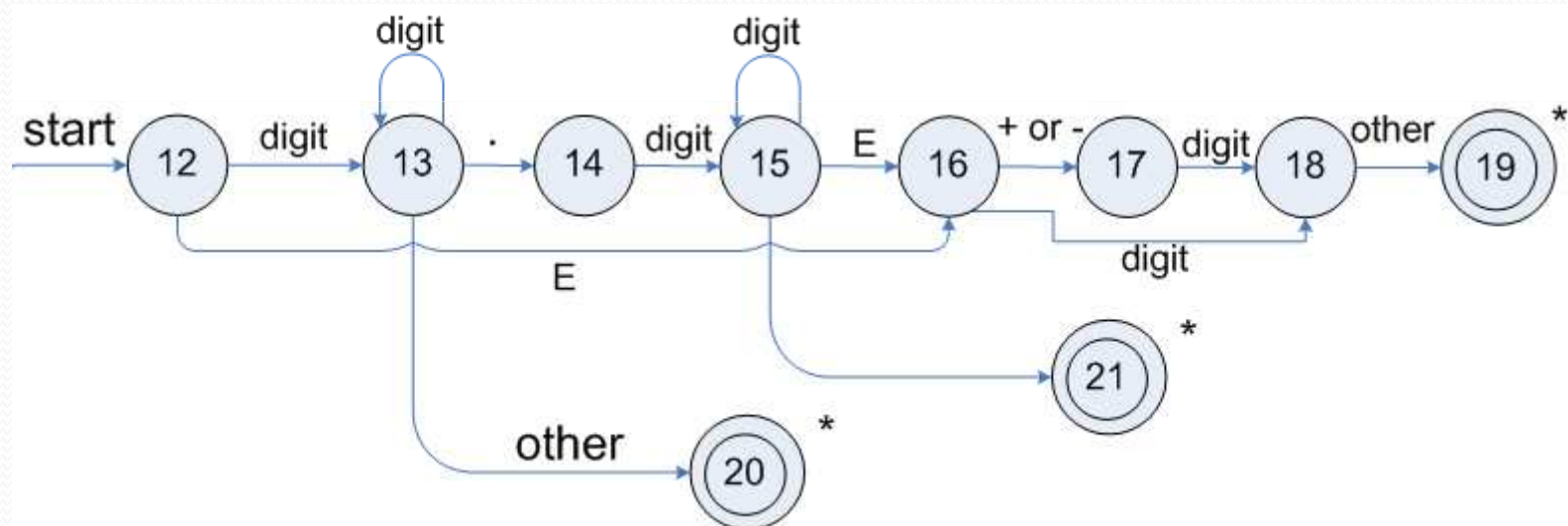
# Transition diagrams (cont...)

- Transition diagram for reserved words and identifiers



# Transition diagrams (cont.)

- Transition diagram for unsigned numbers



# Transition diagrams (cont.)

- Transition diagram for whitespace

