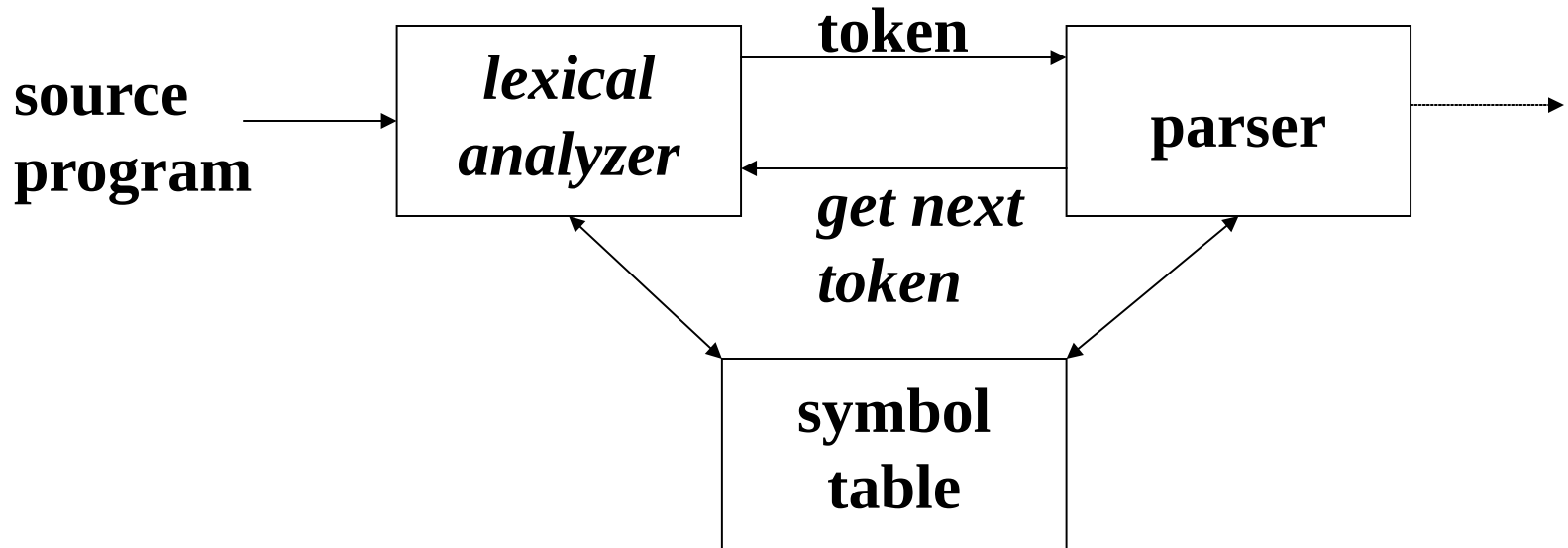


# Lexical Analyzer in Perspective



## Important Issue:

**What are Responsibilities of each Box ?**

**Focus on Lexical Analyzer and Parser**

# Why to separate Lexical analysis and parsing

- o Simplicity of design
- o Improving compiler efficiency
- o Enhancing compiler portability

# Tokens, Patterns, and Lexemes

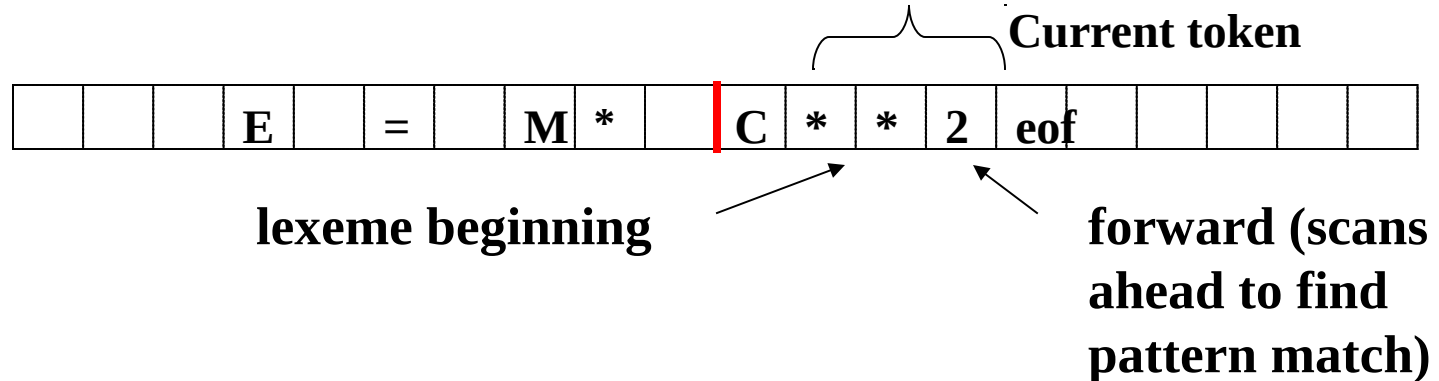
---

- A **token** is a pair a token name and an optional token attribute
- 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

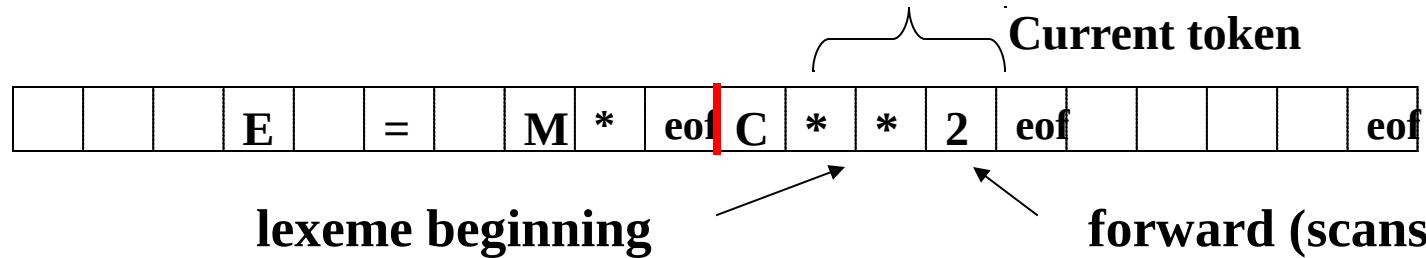
Token	Informal description	Sample lexemes
<b>if</b>	Characters i, f	if
<b>else</b>	Characters e, l, s, e	else
<b>relation</b>	< or > or <= or >= or == or !=	<=, !=
<b>id</b>	Letter followed by letter and digits	pi, score, D2
<b>number</b>	Any numeric constant	3.14159, 0, 6.02e23
<b>literal</b>	Anything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

# Using Buffer to Enhance Efficiency



```
if forward at end of first half then begin
    reload second half ; ← Block I/O
    forward := forward + 1
end
else if forward at end of second half then begin
    reload first half ; ← Block I/O
    move forward to beginning of first half
end
else forward := forward + 1 ;
```

# Algorithm: Buffered I/O with Sentinels



```

forward := forward + 1 ;
if forward is at eof then begin
    if forward at end of first half then begin
        reload second half ; ← Block I/O
        forward := forward + 1
    end
    else if forward at end of second half then begin
        reload first half ; ← Block I/O
        move forward to beginning of first half
    end
    else /* eof within buffer signifying end of input */
        terminate lexical analysis
        2nd eof ⇒ no more input !
end
    
```

forward (scans  
ahead to find  
pattern match)

# Chomsky Hierarchy

---

0 Unrestricted

$$\alpha A \beta \rightarrow \alpha \gamma \beta$$

1 Context-Sensitive

$$| \text{LHS} | \leq | \text{RHS} |$$

2 Context-Free

$$| \text{LHS} | = 1$$

3 Regular

$$| \text{RHS} | = 1 \text{ or } 2 , \\ A \rightarrow a \mid aB, \text{ or} \\ A \rightarrow a \mid Ba$$

# Formal Language Operations

OPERATION	DEFINITION
<i>union</i> of L and M written $L \cup M$	$L \cup M = \{s \mid s \text{ is in } L \text{ or } s \text{ is in } M\}$
<i>concatenation</i> of L and M written LM	$LM = \{st \mid s \text{ is in } L \text{ and } t \text{ is in } M\}$
<i>Kleene closure</i> of L written $L^*$	$L^* = \bigcup_{i=0}^{\infty} L^i$ <p><math>L^*</math> denotes “zero or more concatenations of “ L</p>
<i>positive closure</i> of L written $L^+$	$L^+ = \bigcup_{i=1}^{\infty} L^i$ <p><math>L^+</math> denotes “one or more concatenations of “ L</p>



# Formal Language Operations

## Examples

---

$$L = \{A, B, C, D\} \quad D = \{1, 2, 3\}$$

$$L \cup D = \{A, B, C, D, 1, 2, 3\}$$

$$LD = \{A1, A2, A3, B1, B2, B3, C1, C2, C3, D1, D2, D3\}$$

$$L^2 = \{AA, AB, AC, AD, BA, BB, BC, BD, CA, \dots DD\}$$

$$L^4 = L^2 L^2 = ??$$

$$L^* = \{ \text{All possible strings of } L \text{ plus } \epsilon \}$$

$$L^+ = L^* - \epsilon$$

$$L(L \cup D) = ??$$

$$L(L \cup D)^* = ??$$

# Language & Regular Expressions

---


- A **Regular Expression** is a Set of Rules / Techniques for Constructing Sequences of Symbols (Strings) From an Alphabet.
- Let  $\Sigma$  Be an Alphabet,  $r$  a Regular Expression  
Then  $L(r)$  is the Language That is Characterized by the Rules of  $r$

# Rules for Specifying Regular Expressions:

**fix alphabet  $\Sigma$**

**$\forall \epsilon$  is a regular expression denoting  $\{\epsilon\}$**

- If  $a$  is in  $\Sigma$ ,  $a$  is a regular expression that denotes  $\{a\}$**
- Let  $r$  and  $s$  be regular expressions with languages  $L(r)$  and  $L(s)$ . Then**

- p  
r  
e  
c  
e  
d  
e  
n  
c  
e** 
- (a)  $(r) \mid (s)$  is a regular expression  $\Rightarrow L(r) \cup L(s)$**
  - (b)  $(r)(s)$  is a regular expression  $\Rightarrow L(r) L(s)$**
  - (c)  $(r)^*$  is a regular expression  $\Rightarrow (L(r))^*$**
  - (d)  $(r)$  is a regular expression  $\Rightarrow L(r)$**

**All are Left-Associative. Parentheses are dropped as allowed by precedence rules.**

# EXAMPLES of Regular Expressions

---

$$L = \{A, B, C, D\}$$

$$D = \{1, 2, 3\}$$

$$A \mid B \mid C \mid D = L$$

$$(A \mid B \mid C \mid D)(A \mid B \mid C \mid D) = L^2$$

$$(A \mid B \mid C \mid D)^* = L^*$$

$$(A \mid B \mid C \mid D)((A \mid B \mid C \mid D) \mid (1 \mid 2 \mid 3)) = L(L \cup D)$$

# Algebraic Properties of Regular Expressions

AXIOM	DESCRIPTION
$r \mid s = s \mid r$	$\mid$ is commutative
$r \mid (s \mid t) = (r \mid s) \mid t$	$\mid$ is associative
$(r \ s) t = r (s \ t)$	concatenation is associative
$r (s \mid t) = r s \mid r t$ $(s \mid t) r = s r \mid t r$	concatenation distributes over $\mid$
$\epsilon r = r$ $r \epsilon = r$	$\epsilon$ Is the identity element for concatenation
$r^* = (r \mid \epsilon)^*$	relation between $*$ and $\epsilon$
$r^{**} = r^*$	$*$ is idempotent

# Token Recognition

**How can we use concepts developed so far to assist in recognizing tokens of a source language ?**

**Assume Following Tokens:**

if, then, else, relop, id, num

**Given Tokens, What are Patterns ?**

if ← if

then ← then

else ← else

relop ← < | <= | > | >= | = | <>

id ← letter ( letter | digit )\*

num ← digit <sup>+</sup> ( . digit <sup>+</sup> ) ? ( E ( + | - ) ? digit <sup>+</sup> ) ?

**Grammar:**

*stmt* → |if *expr* then *stmt*  
          |if *expr* then *stmt* else *stmt*  
          |ε

*expr* → *term* relop *term* | *term*

*term* → id | num

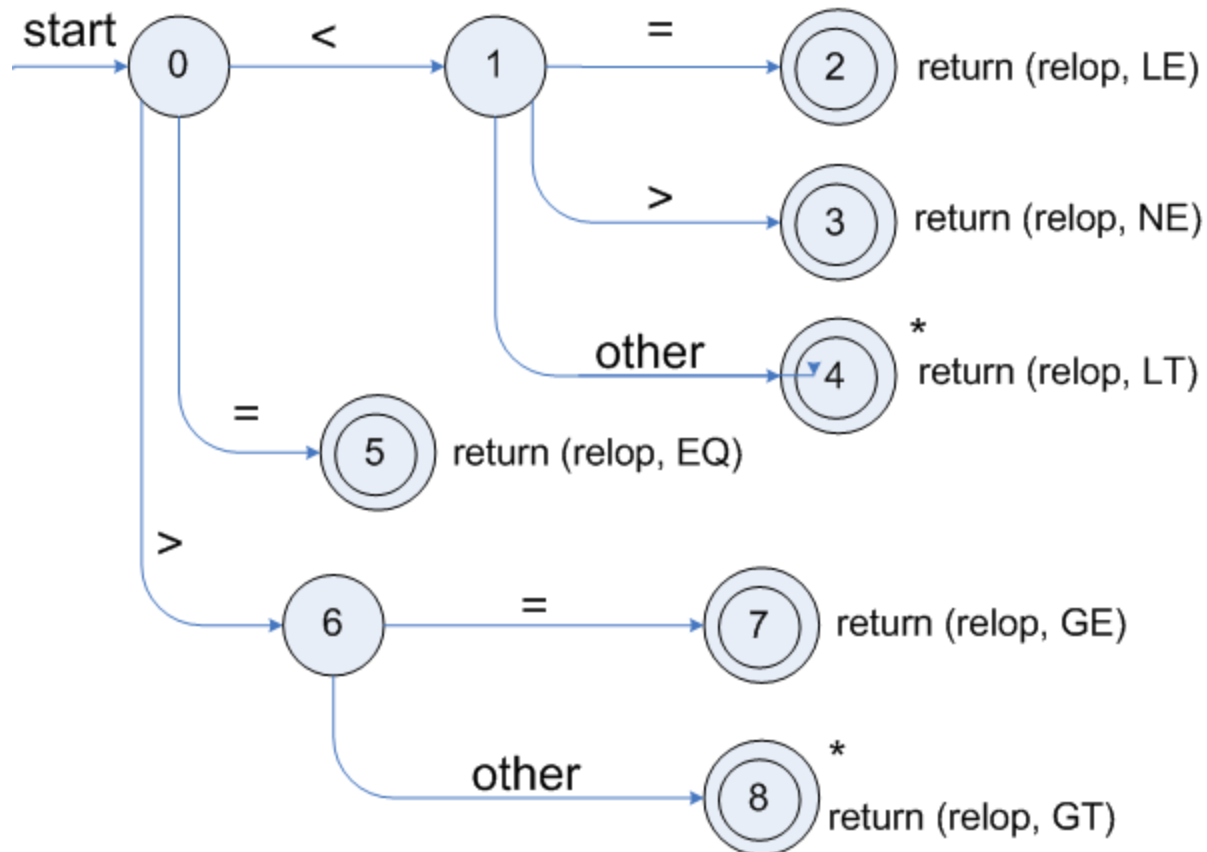
# Overall

Regular Expression	Token	Attribute-Value
WS	-	-
<b>if</b>	<b>if</b>	-
<b>then</b>	<b>then</b>	-
<b>else</b>	<b>else</b>	-
<b>id</b>	<b>id</b>	pointer to table entry
<b>num</b>	<b>num</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>	GE

**Note:** Each token has a unique token identifier to define category of lexemes

# Transition diagrams

## ○ Transition diagram for relop





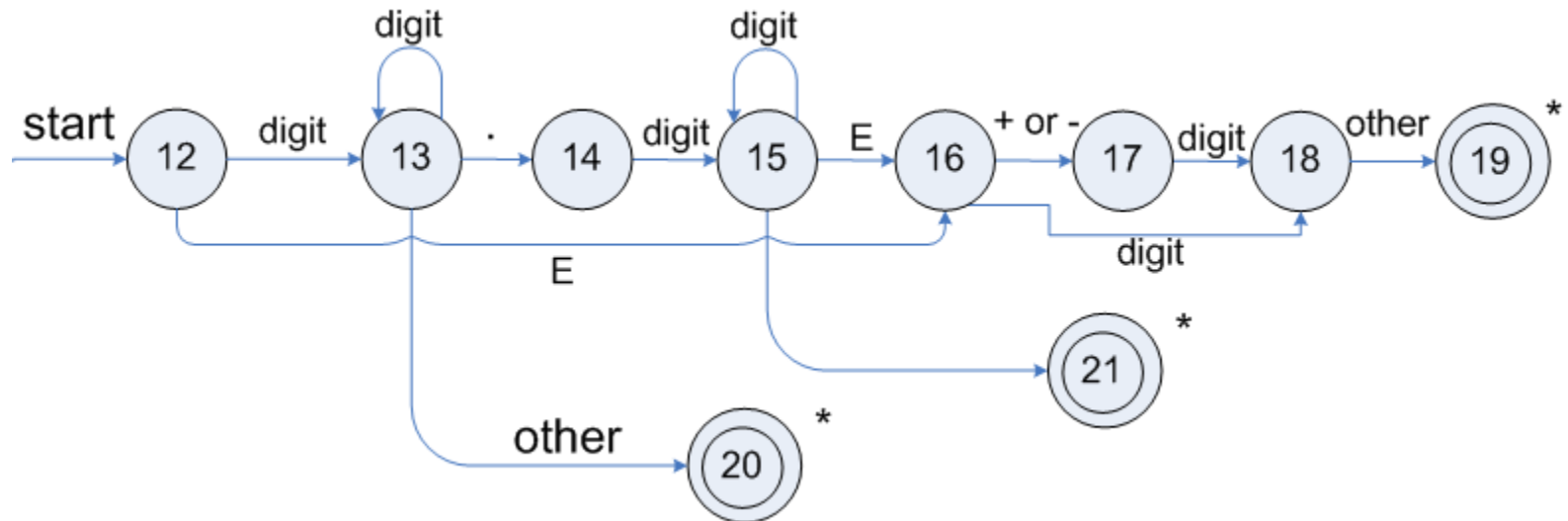
## Transition diagrams (cont.)

- Transition diagram for reserved words and identifiers



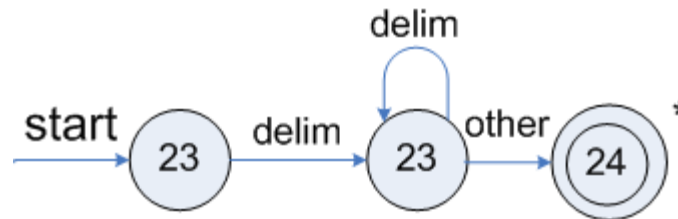
# Transition diagrams (cont.)

- Transition diagram for unsigned numbers

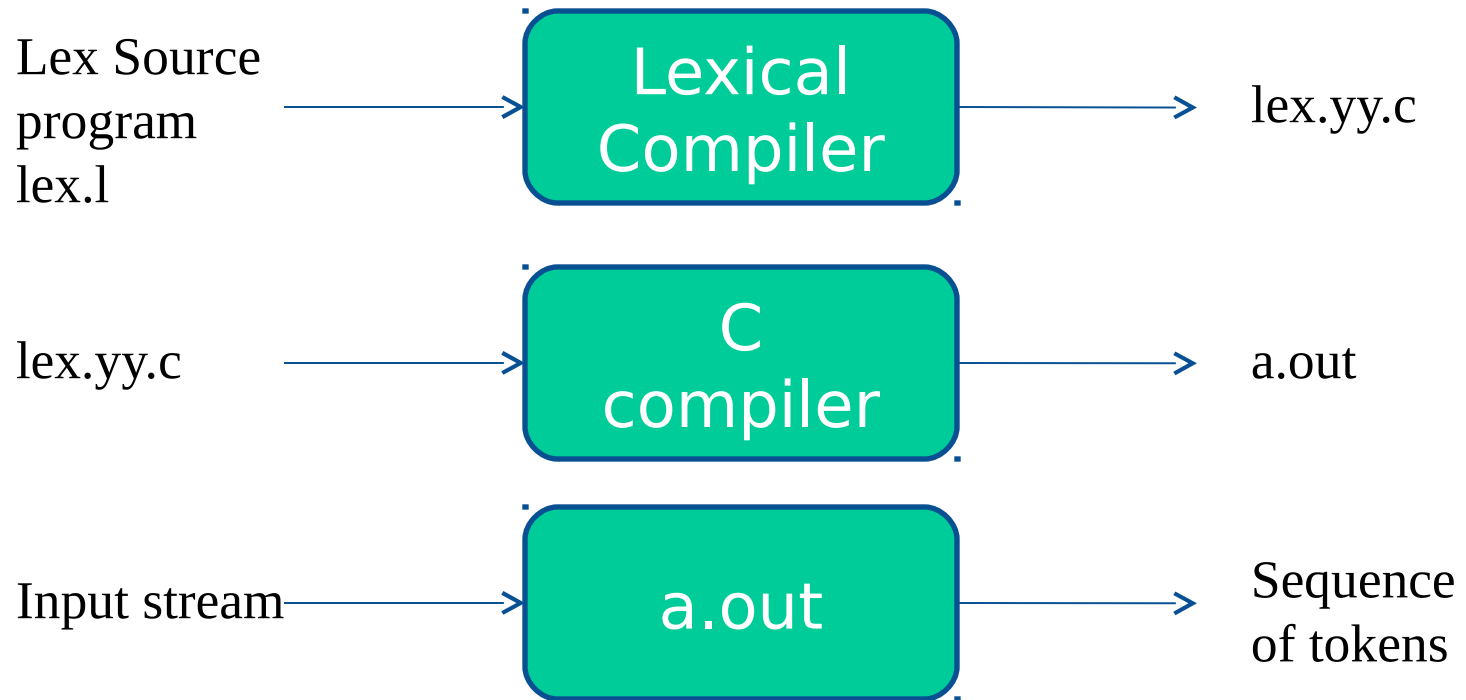


## Transition diagrams (cont.)

- Transition diagram for whitespace



# Lexical Analyzer Generator - Lex



# Lexical errors

- Some errors are out of power of lexical analyzer to recognize:
  - `fi (a == f(x)) ...`
- 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
- Minimal Distance