

# Lexical Analysis

Sukree Sinthupinyo<sup>1</sup>

<sup>1</sup>Department of Computer Engineering  
Chulalongkorn University

14 July 2012

# Outline

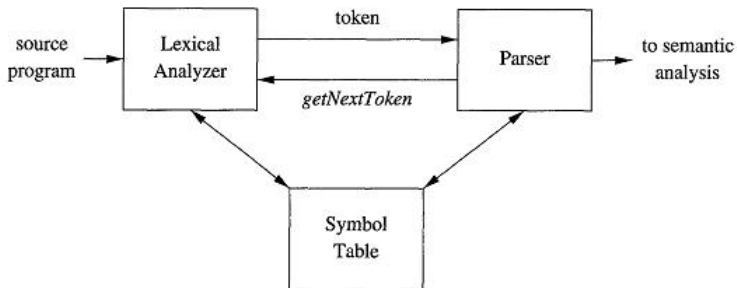
- 1 Introduction
- 2 The Role of the Lexical Analyzer
- 3 Specification of Tokens
  - Regular Expressions
- 4 Recognition of Tokens
  - Transition Diagrams

# Learning Objectives

- Understand definition of *lexeme*, *token*, *etc.*
- Know a method which transforms string into token
- Know syntax of regular expression
- Know concept of transition diagram and code implemented from the diagram

# First step

- The main task is to read the input characters of the source program and export a sequence of tokens.
- It also interacts with the symbol as well.



# First step

- The lexical analyzer must
  - Strip out comments and whitespace.
  - Correlate error messages generated by the compiler with the source program

## Tokens, Patterns, and Lexemes

- 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.
- A *pattern* is a description of the form that the lexemes of a token may take. For the keyword, 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.
- A *lexeme* is a sequence of characters in the source program that matches the pattern for a token.

# Tokens, Patterns, and Lexemes

```
printf("Total = %d\n", score);
```

- `printf` and `score` are lexemes matching the pattern for token **id**
- `"Total = %d\n"` is a lexeme matching **literal**

## Examples of tokens

Token	Informal Description	Sample Lexemes
<b>if</b>	characters <code>i, f</code>	<code>if</code>
<b>else</b>	characters <code>e, l, s, e</code>	<code>else</code>
<b>comparison</b>	<code>&lt;</code> or <code>&gt;</code> or <code>&lt;=</code> or <code>&gt;=</code> or <code>==</code> or <code>!=</code>	<code>&lt;=, !=</code>
<b>id</b>	letter followed by letters and digits	<code>pi, score, D2</code>
<b>number</b>	any numeric constant	<code>3.14, 6.02e23</code>
<b>literal</b>	anything but <code>"</code> , surrounded by <code>"</code> 's	<code>"core"</code>



# General concept of tokens in many programming language

- One token for each keyword. The pattern for a keyword is the same as the keyword itself.
- Tokens for the operators
- One token representing all identifiers
- One or more tokens representing constants, such as numbers and literal strings.
- Tokens for each punctuation symbol, such as left and right parentheses, comma, and semi colon.

## Attributes for Tokens

- Token must have an attribute associated with.
- For example, an **id** must associate with information about identifier; e.g., its lexeme, its type, and the location at which it is first found, is kept in the symbol table.

## An Example of Attributes for Tokens

$E = M * C ** 2$

- **<id, pointer to symbol-table entry for E>**
- **<assign\_op>**
- **<id, pointer to symbol-table entry for M>**
- **<mult\_op>**
- **<id, pointer to symbol-table entry for C>**
- **<exp\_op>**
- **<number, integer value 2 >**

# String and Language

- A *string* over an alphabet is a finite sequence of symbols drawn from that alphabet. The length of string  $s$  is usually written  $|s|$ . The *empty string* is denoted  $\epsilon$ .
- A *language* is any countable set of strings over some fixed alphabet.
- *Concatenation* of string  $x$  and  $y$  is the string formed by appending  $y$  to  $x$ . For example, if  $x = \text{dog}$  and  $y = \text{house}$ , then  $xy = \text{doghouse}$ .
- If we think of concatenation as a product, we can define the "exponentiation" of strings as follows. Define  $s^0$  to be  $\epsilon$ , and for all  $i > 0$ , define  $s^i$  to be  $s^{i-1}s$ . Since  $\epsilon s = s$ , it follows that  $s^i = s$ . Then  $s^2 = ss, s^3 = sss$ , and so on.

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

# Example

- Let  $L$  be the set of letters  $A, B, \dots, Z, a, b, \dots, z$ .
- $D$  be the set of digits  $0, 1, \dots, 9$ .
  - $L \cup D$  is the set of letters and digits with 62 strings of length one.
  - $LD$  is the set of 520 strings of length two.
  - $L^4$  is the set of all 4-letter strings.
  - $L^*$  is the set of all strings of letter, including  $\epsilon$ .
  - $L(L \cup D)^*$  is the set of all strings of letters and digits beginning with a letter.
  - $D^+$  is the set of all strings of one or more digits.

# Outline

- 1 Introduction
- 2 The Role of the Lexical Analyzer
- 3 Specification of Tokens**
  - Regular Expressions
- 4 Recognition of Tokens
  - Transition Diagrams

# Regular Expressions

- If we want to describe the set of valid  $C$  identifiers, we can use the language  $L(L \cup D)$  with the underscore included among the letters.
- If *letter\_* denotes any letter of the underscore, and *digit* stands for any digit, then we could describe the language of  $C$  identifiers by:

$$\textit{letter\_}(\textit{letter\_}|digit)^*$$

- where  $|$  denotes union, the parentheses are used to group subexpressions.



# Regular Expressions

- Language  $L(r)$  is defined recursively from the languages denoted by  $r$ 's subexpressions using alphabet set  $\Sigma$ .
- **BASIS:** There are two rules that form the basis:
  - 1  $\epsilon$  is a regular expression, and  $L(\epsilon)$  is  $\{\epsilon\}$ , that is, the language whose sole member is the empty string.
  - 2 If  $a$  is a symbol in  $\Sigma$ , the  $\mathbf{a}$  is a regular expression, and  $L(\mathbf{a}) = \{a\}$ , that is, the language with one string, of length one, with  $a$  in its one position.

# Regular Expressions

- **INDUCTION:** There are four parts to the induction whereby larger expressions are built from the smaller one. Suppose  $r$  and  $s$  are regular expressions denoting languages  $L(r)$  and  $L(s)$ , respectively.
  - 1  $(r)|(s)$  denotes  $L(r) \cup L(s)$ .
  - 2  $(r)(s)$  denotes  $L(r)L(s)$ .
  - 3  $(r)^*$  denotes  $L(r)^*$ .
  - 4  $(r)$  denotes  $L(r)$ .
- The precedence of operator is  $*$ , concatenation, and  $|$ .
- So  $(a)|((b)^*(c))$  can be written as  $a|b^*c$

# Regular Expressions

- Example

- Let  $\Sigma = \{a, b\}$
- $\mathbf{a|b}$  denotes the language  $\{a, b\}$
- $\mathbf{(a|b)(a|b)}$  denotes  $\{aa, ab, ba, bb\}$
- $\mathbf{a^*}$  denotes  $\{a, aa, aaa, \dots\}$ .
- $\mathbf{(a|b)^*}$  denotes  $\{\epsilon, a, b, aa, ab, ba, bb, aaa, \dots\}$
- $\mathbf{a|a^*b}$  denotes  $\{a, b, ab, aab, aaab, \dots\}$

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

# Definitions

*Regular definition* is a sequence of the form

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

...

$$d_n \rightarrow r_n$$

# Regular Definition Example

- C identifiers are strings of letters, digits, and underscore.

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

## Extensions of Regular Expressions

- $+$ : *One or more instances*
- $?$ : *Zero or one instances*
- $[a_1 a_2 \dots a_n]$ :  $a_1 | a_2 | \dots | a_n$  or  $a_1 - a_n$

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

$digit \rightarrow [0 - 9]$

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

# Example

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

*expr* → *term* **relop** *term*  
      | *term*

*term* → **id**  
      | **number**



# Example

*digit* → [0-9]  
*digits* → *digit*<sup>+</sup>  
*number* → *digits* ( . *digits* ) ? ( E [ + - ] ? *digits* ) ?  
*letter* → [A-Za-z]  
*id* → *letter* ( *letter* | *digit* ) \*  
*if* → **if**  
*then* → **then**  
*else* → **else**  
*relop* → < | > | <= | >= | = | <>

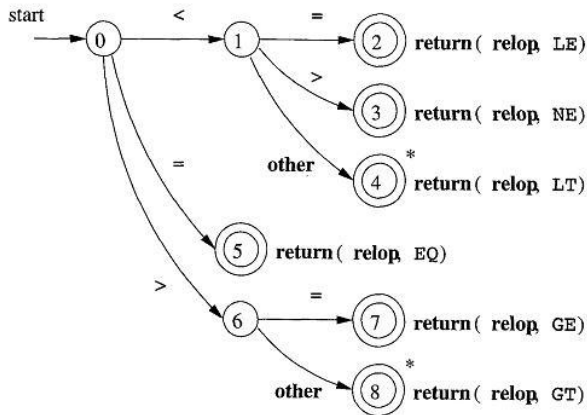
# Tokens, Patterns, and Attribute Values

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>	GE

# Outline

- 1 Introduction
- 2 The Role of the Lexical Analyzer
- 3 Specification of Tokens
  - Regular Expressions
- 4 Recognition of Tokens**
  - Transition Diagrams**

# Transition Diagram for **relop**



# Example Code for **relop**

```
TOKEN getRelop()
{
    TOKEN retToken = new(RELOP);
    while(1) { /* repeat character processing until a return
                or failure occurs */
        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);
        }
    }
}
```