# CS 4300: Compiler Theory

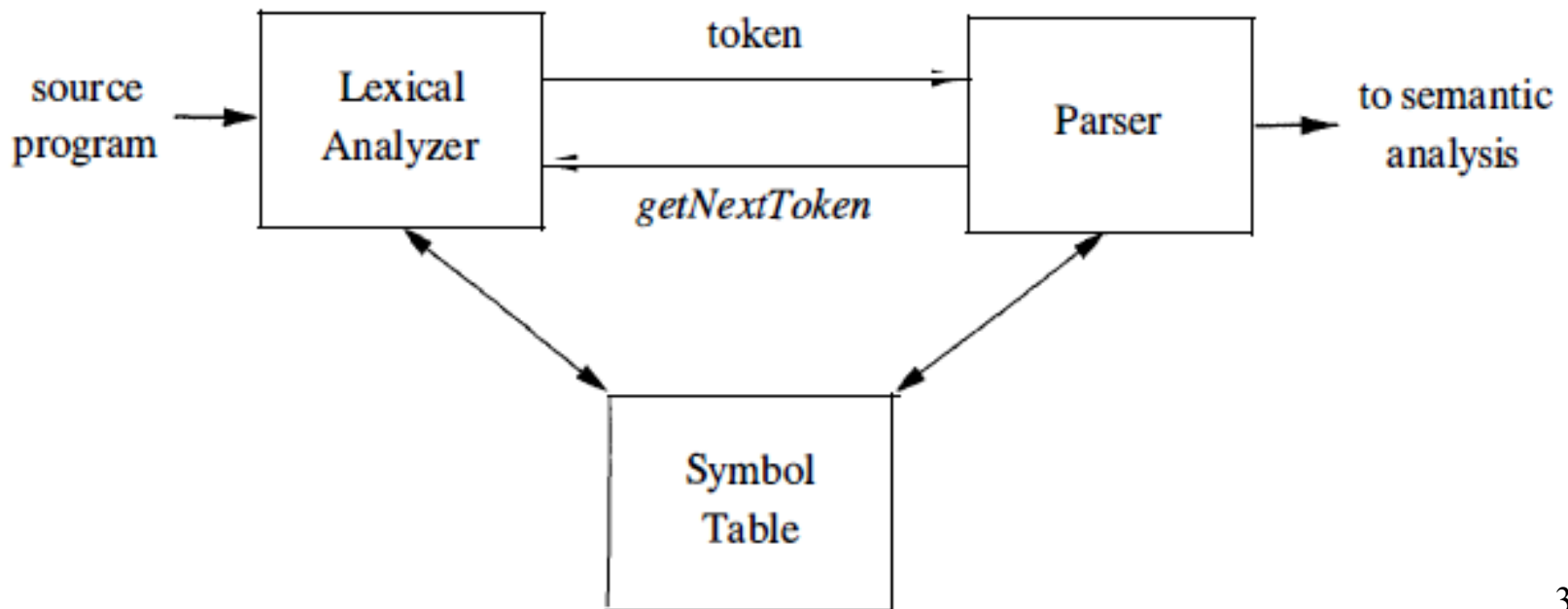# Chapter 3
# Lexical Analysis

*Dr. Xuejun Liang*

# Outlines (Sections)

1. The Role of the Lexical Analyzer
2. Input Buffering (Omit)
3. Specification of Tokens
4. Recognition of Tokens
5. The Lexical -Analyzer Generator Lex
6. Finite Automata
7. From Regular Expressions to Automata
8. Design of a Lexical-Analyzer Generator
9. Optimization of DFA-Based Pattern Matchers*

# 1. The Role of the Lexical Analyzer

- As the first phase of a compiler, the main task of the lexical analyzer is to read the input characters of the source program, group them into lexemes, and produce as output a sequence of tokens for each lexeme in the source program.

source program → Lexical Analyzer —token→ Parser → to semantic analysis

getNextToken (Parser → Lexical Analyzer)

Lexical Analyzer ↔ Symbol Table ↔ Parser

3

# Why Lexical Analysis and Parsing (Syntax Analysis) are Separate

- Simplifies the design of the compiler
  - LL(1) or LR(1) parsing with 1 token lookahead would not be possible (multiple characters/tokens to match)

- Provides efficient implementation
  - Systematic techniques to implement lexical analyzers by hand or automatically from specifications
  - Stream buffering methods to scan input

- Improves portability
  - Non-standard symbols and alternate character encodings can be normalized (e.g. UTF8, trigraphs)

# 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
  - For example: **id** and **num**
- *Lexemes* are the specific character strings that make up a token
  - For example: **abc** and **123**
- *Patterns* are rules describing the set of lexemes belonging to a token
  - For example: "*letter followed by letters and digits*" and "*non-empty sequence of digits*"

# Examples of Tokens

| TOKEN | INFORMAL DESCRIPTION | SAMPLE LEXEMES |
|---|---|---|
| if | characters i, f | if |
| else | characters e, l, s, e | else |
| comparison | < or > or <= or >= or == or != | <=, != |
| id | letter followed by letters and digits | pi, score, D2 |
| number | any numeric constant | 3.14159, 0, 6.02e23 |
| literal | anything but ", surrounded by "'s | "core dumped" |

Token Classes:
1. One token for each keyword
2. Tokens for the operators
3. One token representing all identifiers
4. One or more tokens representing constants
5. Tokens for each punctuation symbol

# Attributes for Tokens

- When more than one lexeme can match a pattern, the lexical analyzer must provide the subsequent compiler phases additional information about the particular lexeme that matched.
    - **number**
- We shall assume that tokens have at most one associated attribute, although this attribute may have a structure that combines several pieces of information
    - E.g. **id** has its lexeme, its type, and the location at which it is first found
    - So the appropriate attribute value for an **id** is a pointer to the symbol-table entry for that identifier (lexeme)

# Example of Attributes for Tokens

- Example: lexemes, token names and associated attribute values for the Fortran statement.

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>

# 3. Specification of Patterns for Tokens: *Definitions*

- An *alphabet* $\Sigma$ is a finite set of symbols (characters)
- A *string s* is a finite sequence of symbols from $\Sigma$
  - $\lvert s \rvert$ denotes the length of string *s*
  - $\varepsilon$ denotes the empty string, thus $\lvert \varepsilon \rvert = 0$
- A *language* is a specific set of strings over some fixed alphabet $\Sigma$

# String Operations

- The *concatenation* of two strings *x* and *y* is denoted by *xy*

- The *exponentation* of a string *s* is defined by

$$s^0 = \varepsilon$$
$$s^i = s^{i-1}s \quad \text{for } i > 0$$

note that $s\varepsilon = \varepsilon s = s$

# Language Operations

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

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

- *Exponentiation*
  $$L^0 = \{\varepsilon\}; \quad L^i = L^{i-1}L$$

- *Kleene closure*
  $$L^* = \cup_{i=0,\ldots,\infty} L^i$$

- *Positive closure*
  $$L^+ = \cup_{i=1,\ldots,\infty} L^i$$

**Example:**
Compute
$L \cup D$
$LD$
$D^4$
$D^*$
$L(L \cup D)^*$
$D^+$

where
$L = \{A, B, ..., Z, a, b, ..., z\}$
and $D = \{0, 1, \ldots 9\}$

# Regular Expressions Over Some Alphabet $\Sigma$

- Basis symbols:
  - $\varepsilon$ is a regular expression denoting language $\{\varepsilon\}$
  - $a \in \Sigma$ is a regular expression denoting $\{a\}$
- If $r$ and $s$ are regular expressions denoting languages $L(r)$ and $L(s)$ respectively, then
  - $r \mid s$ is a regular expression denoting $L(r) \cup L(s)$
  - $rs$ is a regular expression denoting $L(r) \, L(s)$
  - $r^*$ is a regular expression denoting $(L(r))^*$
  - $(r)$ is a regular expression denoting $L(r)$
- A language defined by a regular expression is called a *regular set*

# Precedence of regular expression operations

a) The unary operator * has highest precedence and is left associative.
b) Concatenation has second highest precedence and is left associative
c) | has lowest precedence and is left associative

## Algebraic laws for regular expression operations

| 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 |

**Example 3.4** : Let $\Sigma = \{a, b\}$, what are languages denoted by
The following regular expressions:

**a|b, (a|b)(a|b), a\*, (a|b)\*, a|a\*b**

# Regular Definitions Over Some Alphabet $\Sigma$

- Regular definitions introduce a naming convention with name to regular expression bindings:

$$d_1 \rightarrow r_1$$
$$d_2 \rightarrow r_2$$
$$\dots$$
$$d_n \rightarrow r_n$$

where:
  - Each $d_i$ is a new symbol, not in $\Sigma$ and not the same as any other of the d's, and
  - each $r_i$ is a regular expression over
$$\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$$

# Regular Definitions: Examples

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

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

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

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

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

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

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

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

Numbers: 5280,  0.01234,  6.336E4,  or 1.89E-4.

# Regular Definitions: Extensions

- The following shorthands are often used:

| | |
|---|---|
| One or more instances: + | $r^+ = rr^*$ |
| Zero or one instance: ? | $r? = r \mid \varepsilon$ |
| Character classes: | $[\mathbf{a\text{-}z}] = \mathbf{a} \mid \mathbf{b} \mid \mathbf{c} \mid \ldots \mid \mathbf{z}$ |

- Examples:

$$
\begin{aligned}
letter\_ &\rightarrow [\text{A-Za-z}\_] \\
digit &\rightarrow [\text{0-9}] \\
id &\rightarrow letter\_\ (\ letter\_ \mid digit\ )^*
\end{aligned}
$$

$$
\begin{aligned}
digit &\rightarrow [\text{0-9}] \\
digits &\rightarrow digit^+ \\
number &\rightarrow digits\ (.\ digits)?\ (\ \text{E}\ [\text{+-}]?\ digits\ )?
\end{aligned}
$$

# 4. Recognition of Tokens

$$
\begin{aligned}
stmt \quad &\rightarrow \quad \textbf{if } expr \textbf{ then } stmt \\
&\mid \quad \textbf{if } expr \textbf{ then } stmt \textbf{ else } stmt \\
&\mid \quad \epsilon \\
expr \quad &\rightarrow \quad term \textbf{ relop } term \\
&\mid \quad term \\
term \quad &\rightarrow \quad \textbf{id} \\
&\mid \quad \textbf{number}
\end{aligned}
$$

The terminals of the grammar, which are **if**, **then**, **else**, **relop** , **id**, and **number**, are the names of tokens for lexical analyzer.

# Patterns for tokens of Example 3.8

$$digit \rightarrow [0\text{-}9]$$

$$digits \rightarrow digit^+$$

$$number \rightarrow digits\ (.\ digits)?\ (\ \text{E}\ [\text{+-}]?\ digits\ )?$$

$$letter \rightarrow [\text{A-Za-z}]$$

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

$$if \rightarrow \text{if}$$

$$then \rightarrow \text{then}$$

$$else \rightarrow \text{else}$$

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

# Tokens, patterns, and attribute values

| LEXEMES | TOKEN NAME | ATTRIBUTE VALUE |
|---|---|---|
| Any *ws* | – | – |
| if | if | – |
| then | then | – |
| else | else | – |
| Any *id* | id | Pointer to table entry |
| Any *number* | number | Pointer to table entry |
| < | relop | LT |
| <= | relop | LE |
| = | relop | EQ |
| <> | relop | NE |
| > | relop | GT |
| >= | relop | GE |

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

# Transition Diagrams

$$relop \rightarrow < \mid <= \mid <> \mid > \mid >= \mid =$$



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

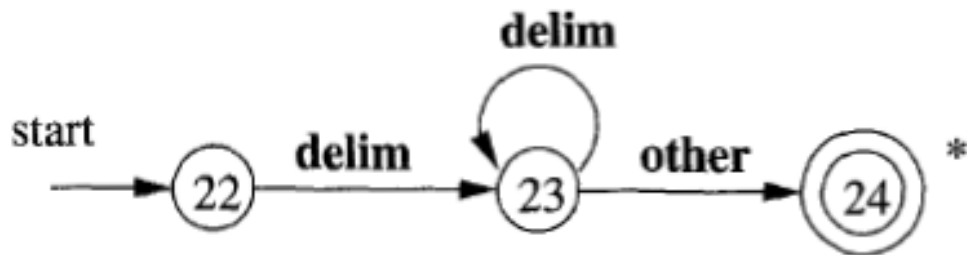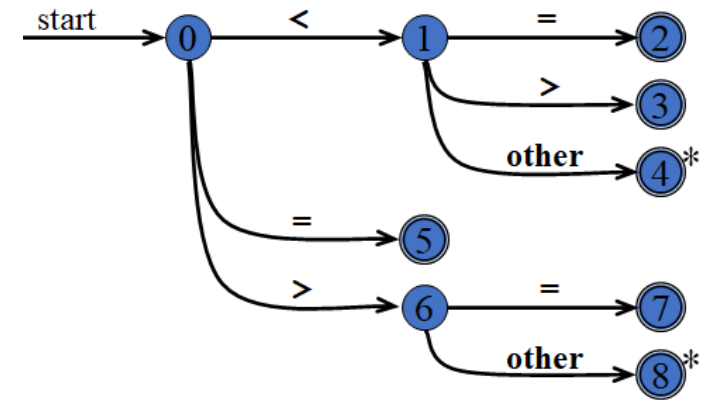# Transition Diagrams (Cont.)

Unsigned number



Whitespace

# Sketch of implementation of relop transition diagram



```
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);
        }
    }
}
```