



中山大學  
SUN YAT-SEN UNIVERSITY



国家超级计算广州中心  
NATIONAL SUPERCOMPUTER CENTER IN GUANGZHOU

# Compilation Principle 编译原理

---

## 第1讲：词法分析(1)

张献伟

[xianweiz.github.io](https://xianweiz.github.io)

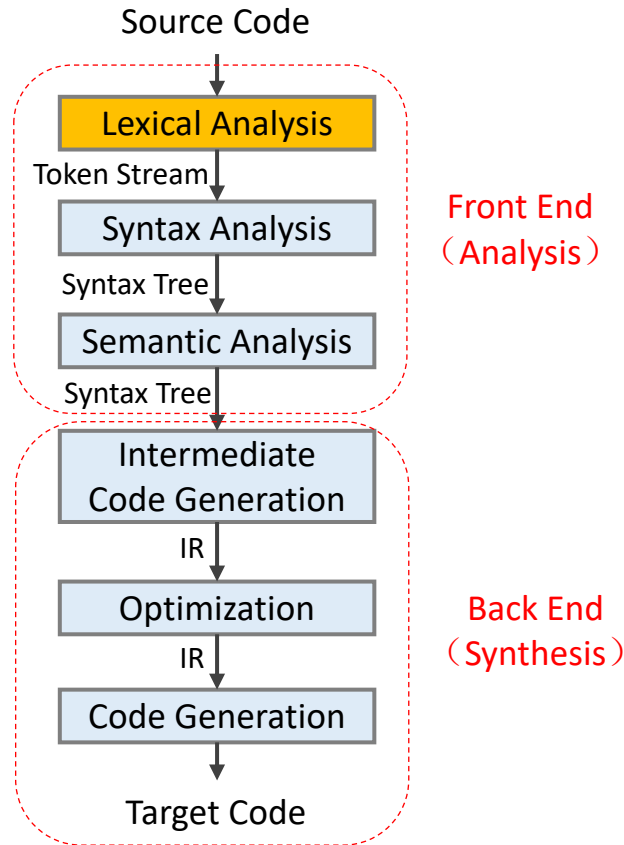
DCS290, 3/2/2021



中山大學  
SUN YAT-SEN UNIVERSITY



# Structure of a Typical Compiler[结构]



# What is Lexical Analysis[词法分析]?

- Example:

```
/* simple example */  
if (i == j)  
    z = 0;  
else  
    z = 1;
```

- Input: a string of characters[输入]
  - “if (i == j) \n \t \tz = 0; \telse\n\tz = 1; \n”
- Goal: partition the string into a set of substrings[目标]
  - Those substrings are **tokens**
- Steps[步骤]
  - Remove comments
  - Identify substrings: ‘if’ ‘(’ ‘i’ ‘==’ ‘j’ .....
  - Identify **token classes**: (keyword, ‘if’), (LPAR, ‘(’), (id, ‘i’) .....

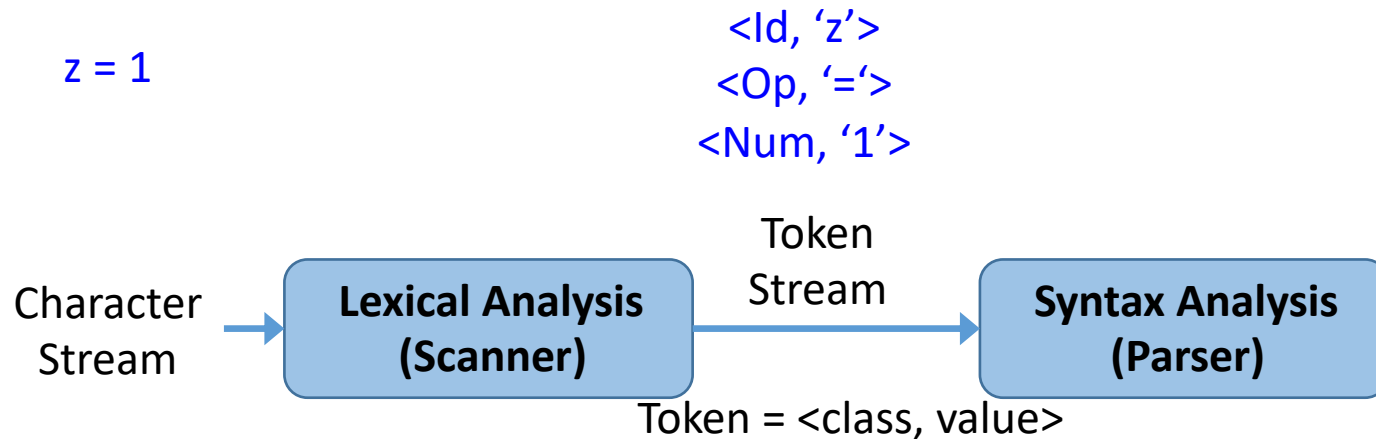
# What is a token[词]?

---

- **Token**: a “word” in language (smallest unit with meaning)
  - Categorized into classes according to its role in language
  - Token classes in English[自然语言]
    - Noun, verb, adjective, ...
  - Token classes in a programming language[编程语言]
    - Number, keyword, whitespace, identifier, ...
- Each **token class** corresponds to a set of strings
  - **Numbers**: a non-empty string of digits
  - **Keyword**: a fixed set of reserved words (“for”, “if”, “else”, ...)
  - **Whitespace**: a non-empty sequence of blanks, tabs, newlines
  - **Identifier**: user-defined name of an entity to identify (Q: what are the rules in C language?)

# Lexical Analysis: Tokenization[分词]?

- Lexical analysis is also called **Tokenization** (also called Scanner)[词法分析也称为扫描器]
  - Partition input string into a sequence of tokens
  - Classify each token according to roles (token class)
    - **Lexeme**: an instance of the corresponding token class, e.g. 'z', '=', '1'
- Pass tokens to syntax analyzer (also called Parser)[分析器]
  - Parser relies on token classes to identify roles (e.g., a keyword is treated differently than an identifier)



# Lexical Analyzer: Design[词法分析器设计]

- Define a finite set of token classes
  - Describe all items of interest
  - Depends on language, design of parser
  - “if (i == j) \n t \tz = 0; \telse\n \tz = 1; \n”
    - Keyword, identifier, whitespace, integer
- Label which string belongs to which token class

```
if (i == j)
    z = 0;
else
    z = 1;
```

'==' or '='?

keyword or identifier?

# Lexical Analyzer: Implementation[实现]

---

- An implementation must do two things
  - Recognize the token class the substring belongs to
  - Return the value or lexeme of the token
- A token is a tuple (class, lexeme)[二元组]
- The lexer usually discards “non-interesting” tokens that don’t contribute to parsing[丢弃无意义词]
  - e.g., whitespace, comments
- If token classes are non-ambiguous, tokens can be recognized in a single left-to-right scan of input string
- Problem can occur when classes are ambiguous[歧义]

# Ambiguous Tokens in C++

- C++ template syntax
  - Foo<Bar>
- C++ stream syntax
  - cin >> var

- Ambiguity

- Foo<Bar<Bar>>>
- cin >>> var
- Q: Is '>>>' a stream operator or two consecutive brackets?

```
Template <typename T>
```

```
T getMax(T x, T y) {  
    return (x > y) ? x : y;  
}
```

```
int main (int argc, char* argv[]) {  
    getMax<int>(3, 7);  
    getMax<double>(3.0, 2.0);  
    getMax<char>('g', 'e');  
  
    return 0;  
}
```



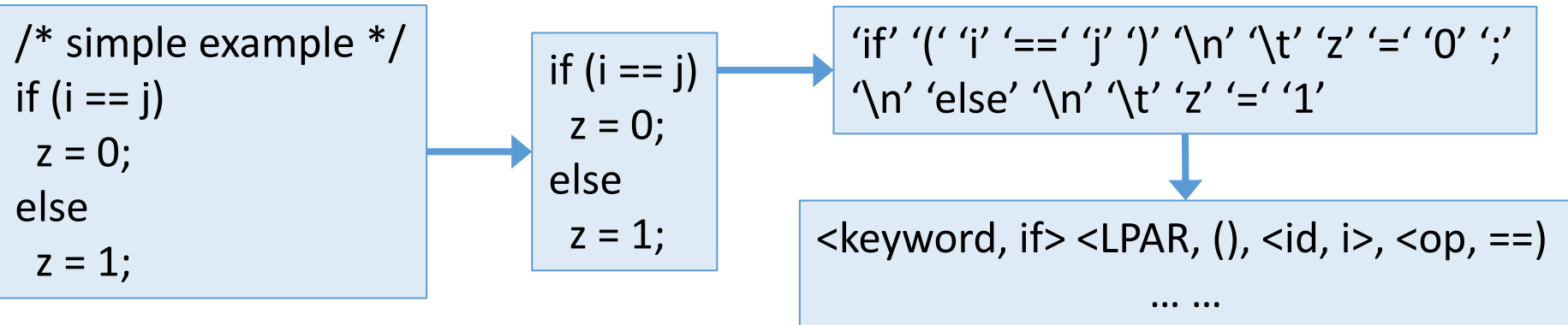
# Look Ahead[展望]

---

- “look ahead” may be required to resolve ambiguity[展望消除歧义]
  - Extracting some tokens requires looking at the larger context or structure
  - Structure emerges only at parsing stage with parse tree
  - Hence, sometimes feedback from parser needed for lexing
    - This complicates the design of lexical analysis
    - Should minimize the amount of look ahead
- Usually tokens do not overlap[通常无重叠]
  - Tokenizing can be done in one pass w/o parser feedback
  - Clean division between lexical and syntax analyses

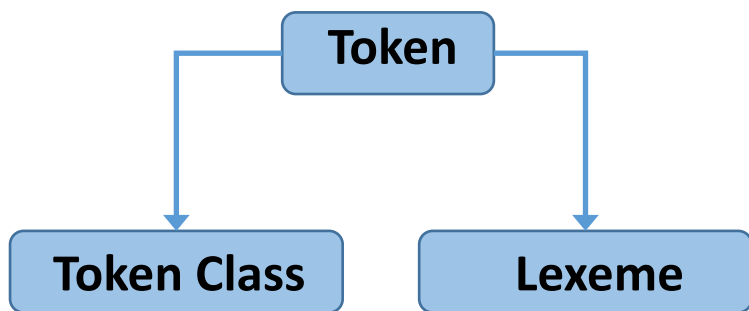
# Summary: Lexer

- Lexical Analysis
  - Partition the input string to lexeme
  - Identify the token class of each lexeme
- Left-to-right scan => look ahead may be required
  - In reality, lookahead is always needed
  - The amount of lookahead should be minimized



# Token Specification[定义]

- Recognizing token class: how to describe string patterns
  - i.e., which set of strings belong to which token class?
  - Use regular expressions [正则表达式] to define token class
- Regular Expression is a good way to specify tokens
  - Simple yet powerful (able to express patterns)
  - Tokenizer implementation can be generated automatically from specification (using a translation tool)
  - Resulting implementation is provably efficient



String patterns  
describing the class

# Language: Definition

---

- **Alphabet**  $\Sigma$  [字母表]: a finite set of symbols
  - Symbol: letter, digit, punctuation, ...
  - Example:  $\{0, 1\}$ ,  $\{a, b, c\}$ , ASCII
- **String** [串]: a finite sequence of symbols drawn from  $\Sigma$ 
  - Example: aab (length = 3),  $\epsilon$  (empty string, length = 0)
- **Language** [语言]: a set of strings of the characters drawn from  $\Sigma$ 
  - $\Sigma = \{0, 1\}$ , then  $\{\}$ ,  $\{01, 10\}$ ,  $\{1, 11, 1111, \dots\}$  are all languages over  $\Sigma$
  - $\{\epsilon\}$  is a language
  - $\Phi$ , empty set is also a language

# Language: Example

---

- Examples:
  - Alphabet  $\Sigma$  = (set of) English characters
  - Language  $L$  = (set of) English sentences
  - Alphabet  $\Sigma$  = (set of) Digits, +, -
  - Language  $L$  = (set of) Integer numbers
- Languages are subsets of all possible strings
  - Not all strings of English characters are sentences
  - Not all sequences of digits and signs are integers

# Regular Expressions & Languages[正则]

---

- Need a notion to specify strings in a particular language
  - More complex languages need more complex notations
- **Regular Expression** is a simple notation
  - Can express simple patterns (e.g., repeating sequences)
  - Not powerful enough to express English (or even C)
  - But powerful enough to express tokens (e.g., identifiers)
- Languages that can be expressed using regular expressions are called **Regular Languages**
- More complex languages and expressions will be covered later

# Atomic REs[原子表达式]

---

- Atomic
  - Smallest RE that cannot be broken down further
- **Epsilon or  $\epsilon$**  character denotes a zero length string
  - $\epsilon = \{""\}$
- **Single character** denotes a set of one string
  - $'c' = \{“c”\}$
- Empty set is  $\{ \} = \phi$ , not the same as  $\epsilon$ 
  - $\text{Size}(\phi) = 0$
  - $\text{Size}(\epsilon) = 1$
  - $\text{Length}(\epsilon) = 0$

# Compound REs[组合表达式]

---

- **Union[并]**: if A and B are REs, then
$$A|B = \{ s \mid s \in A \text{ or } s \in B \}$$
- **Concatenation[连接]** of sets/strings
$$AB = \{ ab \mid a \in A \text{ and } b \in B \}$$
- **Iteration[迭代]** (Kleene closure)
$$A^* = \bigcup_{i \geq 0} A^i \text{ where } A^i = A \dots A \text{ (i times)}$$

in particular

$$A^* = \{\epsilon\} + A + AA + AAA + \dots$$
$$A^+ = A + AA + AAA + \dots = AA^*$$
- **(A)  $\equiv$  A**: A is a RE



# RE and RL

---

- The regular expressions (REs) over  $\Sigma$  are the total set of expressions that can be constructed using components:
  - $\epsilon$
  - $'c'$  where  $c \in \Sigma$
  - $A|B$  where  $A, B$  are REs over  $\Sigma$
  - $AB$  where  $A, B$  are REs over  $\Sigma$
  - $A^*$  where  $A$  is a RE over  $\Sigma$
- The regular languages (RLs) over  $\Sigma$  are the total set of languages that can be expressed using REs:
  - $L(\epsilon) = \{''''\}$
  - $L('c') = \{''c''\}$
  - $L(A|B) = L(A) \cup L(B)$
  - $L(AB) = \{ab \mid a \in L(A) \text{ and } b \in L(B)\}$
  - $L(A^*) = \bigcup_{i \geq 0} L(A^i)$

# Operator Precedence[优先级]

---

- RE operator precedence

- $(A)$
- $A^*$
- $AB$
- $A|B$

- Example:  $ab^*c|d$

- $a(\underline{b^*})c|d$
- $(\underline{a(b^*)})c|d$
- $(\underline{(a(b^*))c})|d$

# Common REs[常用表达]

---

- **At least one:**  $A^+ \equiv AA^*$
- **Union:**  $A|B \equiv A+B$
- **Option:**  $A? \equiv A + \varepsilon$
- **Range:**  $'a' + 'b' + \dots + 'z' \equiv [a-z]$
- **Excluded range:** complement of  $[a-z] \equiv [^a-z]$

# RE Examples

Regular Expression	Explanation
$a^*$	0 or more a's ( $\epsilon$ , a, aa, aaa, aaaa, ...)
$a^+$	1 or more a's (a, aa, aaa, aaaa, ...)
$(a b)(a b)$	(aa, ab, ba, bb)
$(a b)^*$	all strings of a's and b's (including $\epsilon$ )
$(aa ab ba bb)^*$	all strings of a's and b's of even length
$[a-zA-Z]$	shorthand for " $a b \dots z A B \dots Z$ "
$[0-9]$	shorthand for " $0 1 2 \dots 9$ "
$0([0-9])^*0$	numbers that start and end with 0
$1^*(0 \epsilon)1^*$	binary strings that contain at most one zero
$(0 1)^*00(0 1)^*$	all binary strings that contain '00' as substring

- Q: are  $(a|b)^*$  and  $(a^*b^*)^*$  equivalent?

# More Examples

---

- Keywords: 'if' or 'else' or 'then' or 'for' ...
  - RE = 'i''f' + 'e''l''s''e' + ... = 'if' + 'else' + 'then' + ...
- Numbers: a non-empty string of digits
  - digit = '0' + '1' + '2' + '3' + '4' + '5' + '6' + '7' + '8' + '9'
  - integer = digit digit\*
  - Q: is '000' an integer?
- Identifier: strings of letters or digits, starting with a letter
  - letter = 'a' + 'b' + ... 'z' + 'A' + 'B' + ... + 'Z' = [a-zA-Z]
  - RE = letter(letter + digit)\*
  - Q: is the RE valid for identifiers in C?
- Whitespace: a non-empty sequence of blanks, newline and tabs
  - (' ' + '\n' + '\t')+

# REs in Programming Language

Symbol	Meaning		
<code>\d</code>	Any decimal digit, i.e. [0-9]		
<code>\D</code>	Any non-digit char, i.e., [^0-9]		
<code>\s</code>	Any whitespace char, i.e., [ \t\n\r\f\v]		
<code>\S</code>	Any non-whitespace char, i.e., [^ \t\n\r\f\v]		
<code>\w</code>	Any alphanumeric char, i.e., [a-zA-Z0-9_]		
<code>\W</code>	Any non-alphanumeric char, i.e., [^a-zA-Z0-9_]		
<code>.</code>	Any char	<code>\.</code>	Matching “.”
<code>[a-f]</code>	Char range	<code>[^a-f]</code>	Exclude range
<code>^</code>	Matching string start	<code>\$</code>	Matching string end
<code>(...)</code>	Capture matches		

<https://docs.python.org/3/howto/regex.html>

# Lexical Specification of a Language

---

- **S0**: write a regex for the lexemes of each token class
  - Numbers =  $\text{digit}^+$
  - Keywords = 'if' + 'else' + ...
  - Identifiers =  $\text{letter}(\text{letter} + \text{digit})^*$
- **S1**: construct  $R$ , matching all lexemes for all tokens
  - $R = \text{numbers} + \text{keywords} + \text{identifiers} + \dots = R_1 + R_2 + R_3 + \dots$
- **S2**: let input be  $x_1 \dots x_n$ , for  $1 \leq i \leq n$ , check  $x_1 \dots x_i \in L(R)$
- **S3**: if successful, then we know  $x_1 \dots x_i \in L(R_j)$  for some  $j$
- **S4**: remove  $x_1 \dots x_i$  from input and go to step S2

# Lexical Specification of a Language

---

- How much input is used?
  - $x_1 \dots x_i \in L(R)$ ,  $x_1 \dots x_j \in L(R)$ ,  $i \neq j$
  - Which one do we want? (e.g., '==' or '=')
  - Maximal match: always choose the longer one[最长匹配]
- Which token is used if more than one matches?
  - $x_1 \dots x_i \in L(R)$  where  $R = R_1 + R_2 + \dots + R_n$
  - $x_1 \dots x_i \in L(R_m)$ ,  $x_1 \dots x_i \in L(R_n)$ ,  $m \neq n$
  - E.g., keywords = 'if', identifier = letter(letter+digit)\*
  - Keyword has higher priority
  - Rule of thumb: choose the one listed first[次序]
- What if no rule matches?
  - $x_1 \dots x_i \notin L(R) \rightarrow \text{Error}$



# Summary: RE

---

- We have learnt how to specify tokens for lexical analysis[定义token]
  - Regular expressions
  - Concise notations for the string patterns
- Used in lexical analysis with some extensions[适度扩展]
  - To resolve ambiguities
  - To handle errors
- REs is only a language specification[只是定义了语言]
  - An implementation is still needed
  - Next: to construct a token recognizer for languages given by regular expressions – by using **finite automata**