



Compilation Principle 编译原理

第4讲: 词法分析(4)

张献伟

xianweiz.github.io

DCS290, 3/7/2023





Quiz Questions



- Q1: write RE for binary numbers that are multipliers of 4?
 (0|1)*00
- Q2: lexical analysis of 'if (a != b'?

(keyword, 'if'), (sym, '('), (id, 'a'), (sym, '!='), (id, 'b')

- Q3: regard lexer implementation, why NFA → DFA?
 Trade-off space for speed; DFA is more efficient
- Q4: RE of the FA?
 (0|1)*1



• Q5: start state of the equivalent DFA?

	\bigwedge^{1}
1	ABC
0	

 ϵ -closure(A) = {A, B}

 ϵ -closure(move({AB}, 0)) = ϵ -closure({A}) \Longrightarrow {A, B}

 ϵ -closure(move({AB}, 1)) = ϵ -closure({A,C}) \Longrightarrow {A, B, C}



1

ABC

ABC

 $\mathbf{0}$

AB

AB

AB

ABC

Lexical Analyzer Generated by Lex

- A Lex program is turned into a transition table and actions, which are used by a FA simulator
- Automaton recognizes matching any of the patterns





Lex: Example

- Three patterns, three NFAs
- Combine three NFAs into a single NFA
 - Add start state 0 and ε-transitions





Lex: Example (cont.)







Lex: Example (cont.)

- NFA's for lexical analyzer
- Input: aaba
 - $-\epsilon$ -closure(0) = {0, 1, 3, 7}

abb a*b+

а

- Empty states after reading the fourth input symbol
 - There are no transitions out of state 8
 - Back up, looking for a set of states that include an accepting state
- State 8: a*b+ has been matched
 - □ Select aab as the lexeme, execute action₃

\square Return to parser indicating that token w/ pattern p₃=a*b+ has been found



Lex: Example (cont.)

- DFA's for lexical analyzer
- Input: abba
 - Sequence of states entered: 0137 \rightarrow 247 \rightarrow 58 \rightarrow 68 a^{*b+}
 - At the final *a*, there is no transition out of state 68

• 68 itself is an accepting state that reports pattern $p_2 = abb$





а

abb

How Much Should We Match?[匹配多少]

- In general, find the longest match possible
 - We have seen examples
 - One more example: input string aabbb ...
 - Have many prefixes that match the third pattern
 - Continue reading b's until another a is met
 - Report the lexeme to be the intial a's followed by as many b's as there are

 $\{action_1\}$

{ action, }

 $\{action_3\}$

а

abb

a*b+

- If same length, rule appearing first takes precedence
 - String *abb* matches both the second and third
 - We consider it as a lexeme for p_2 , since that pattern listed first

ptn1 ptn2 ptn3	a abb a*b+	ptn1 ptn2 ptn3	a abb a*b+
%%	<ptn2, abb=""></ptn2,>	%%	<ptn3, abb=""></ptn3,>
<pre>{ptn1} {ptn2} {ptn3}</pre>	<pre>{ printf("\n<%s, %s>", "ptn1", yytext); } { printf("\n<%s, %s>", "ptn2", yytext); } { printf("\n<%s, %s>", "ptn3", yytext); }</pre>	<pre>{ptn1} {ptn3} {ptn2}</pre>	<pre>{ printf("\n<%s, %s>", "ptn1", yytext); } { printf("\n<%s, %s>", "ptn3", yytext); } { printf("\n<%s, %s>", "ptn2", yytext); }</pre>
	ヤ山大學 IN YAT-SEN UNIVERSITY	8	

How to Match Keywords?[匹配关键字]

- Example: to recognize the following tokens
 - Identifiers: letter(letter|digit)*
 - Keywords: if, then, else
- Approach 1: make REs for keywords and place them before REs for identifiers so that they will take precedence
 - Will result in more bloated finite state machine
- Approach 2: recognize keywords and identifiers using same RE but differentiate using special keyword table
 - Will result in more streamlined finite state machine
 - But extra table lookup is required
- Usually approach 2 is more efficient than 1, but you can implement approach 1 in your projects for simplicity



The Conversion Flow[转换流程]

- Outline: RE \rightarrow NFA \rightarrow DFA \rightarrow Table-driven Implementation
 - ③ Converting DFAs to table-driven implementations
 - 1 Converting REs to NFAs (M-Y-T algorithm)
 - –
 ⁽²⁾ Converting NFAs to DFAs (subset construction)
 - ⁽³⁾ DFA minimization (partition algorithm)



Beyond Regular Languages

- Regular languages are expressive enough for tokens

 Can express identifiers, strings, comments, etc.
- However, it is the weakest (least expressive) language
 - Many languages are not regular
 - C programming language is not
 - The language matching braces "{{{...}}}" is also not
 - FA cannot count # of times char encountered
 - $\Box L = \{a^n b^n \mid n \ge 1\}$
 - Crucial for analyzing languages with nested structures (e.g. nested for loop in C language)
- We need a more powerful language for parsing
 - Later, we will discuss context-free languages (CFGs)



RE/FA is NOT Powerful Enough

- $L = \{a^nb^n \mid n \ge 1\}$ is **NOT** a Regular Language
 - Suppose L were the language defined by regular expression
 - Then we could construct a DFA D with k states to accept L
 - Since D has only k states, for an input beginning with more than k a's,
 D must enter some state twice, say s_i
 - Suppose that the path from s_i back to itself is labeled with a^{j-i}
 - Since *aⁱbⁱ* is in *L*, there must be a path labeled *bⁱ* from *s_i* to an accepting state *f*
 - But, there is also a path from s_0 through s_i to f labelled $a^i b^i$
 - Thus, *D* also accepts *aⁱbⁱ*, which is not in *L*, contradicting the assumption that *L* is the language accepted by *D*



RE/FA is NOT Powerful Enough(cont.)

- $L = \{a^nb^n \mid n \ge 1\}$ is not a Regular Language
 - Proof → Pumping Lemma (泵引理)
 - FA does not have any memory (FA cannot count)
 - □ The above *L* requires to keep count of a's before seeing b's
- Matching parenthesis is not a RL
- Any language with nested structure is not a RL
 if ... if ... else ... else
- Regular Languages
 - Weakest formal languages that are widely used





Compilation Principle 编译原理

张献伟

<u>xianweiz.github.io</u>

DCS290, 3/7/2023





Compilation Phases[编译阶段]







Example

```
void 'void'
                       void main() {
  • $vim test.c
                                                                  identifier 'main'
                       int;
                                                                  l paren '('
                        int a,;
                                                                  r paren ')'
                                                                  l_brace '{'
                        int b, c;
                                                                  int 'int'
                                                                  semi ';'

    $clang -cc1 -dump-tokens ./test.c

                                                                  int 'int'
                                                                  identifier 'a'

    $clang -o test test.c

                                                                  comma ','
                                                                  semi ';'
test.c:1:1: warning: return type of 'main' is not 'int' [-Wmain-return<mark>int 'int'</mark>
void main() {
                                                                  identifier 'b'
                                                                  comma ','
test.c:1:1: note: change return type to 'int'
void main() {
                                                                  identifier 'c'
Anna
                                                                  semi ';'
int
                                                                  r_brace '}'
test.c:2:3: warning: declaration does not declare anything [-Wmissing-
                                                                  eof ''
ns]
  int;
test.c:3:9: error: expected identifier or '('
  int a,;
2 warnings and 1 error generated.
                                            16
```

Syntax Analysis[语法分析]

- Second phase of compilation[第二阶段]
 - Also called as parser
- Parser obtains a string of tokens from the lexical analyzer[以token作为输入]
 - Lexical analyzer reads the chars of the source program, groups them into lexically meaningful units called lexemes
 - and produces as output tokens representing these lexemes
 - Token: <token name, attribute value>
 - Token names are used by parser for syntax analysis
 □ tokens → parse tree/AST
- Parse tree[分析树]
 - Graphically represent the syntax structure of the token stream



Parsing Example

- Input: if(x==y) ... else ...[源程序输入]
- Parser input (Lexical output)[语法分析输入] KEY(IF) '(' ID(x) OP('==') ID(y) ')' ... KEY(ELSE) ...
- Parser output[语法分析输出]







Parsing Example (cont.)

- Example: <id, x> <op, *> <op, %>
 - Is it a valid token stream in C language? **YES**
 - Is it a valid statement in C language (x *%)? NO
- Not every sequence of tokens are valid
 - Parser must distinguish between valid and invalid token sequence
- We need a method to describe what is valid sequence?
 To specify the syntax of a programming language



How to Specify Syntax?

- How can we specify a syntax with nested structures?
 - Is it possible to use RE/FA?
 - L(Regular Expression) ≡ L(Finite Automata)
- RE/FA is not powerful enough

 $-L = \{a^nb^n \mid n \ge 1\}$ is not a Regular Language

- Example: matching parenthesis: # of '(' == # of ')'
 - $(x+y)*z \qquad \checkmark \\ ((x+y)+y)*z \qquad \checkmark \\ (...(((x+y)+y)+y)...) \qquad \checkmark$







What Language Do We Need?

- C-language syntax: Context Free Language (CFL)[上下文无 关语言] e.g., 'else' is always 'else', wherever you place it
 - A broader category of languages that includes languages with nested structures
- Before discussing CFL, we need to learn a more general way of specifying languages than RE, called Grammars[文 法]
 - Can specify both RL and CFL
 - and more ...
- Everything that can be described by a regular expression can also be described by a grammar
 - Grammars are most useful for describing nested structures





Concepts

- Language[语言]
 - Set of strings over alphabet
 - String: finite sequence of symbols
 - Alphabet: finite set of symbols
- Grammar[文法]
 - To systematically describe the syntax of programming language constructs like expressions and statements
- Syntax[语法]
 - Describes the proper form of the programs
 - Specified by grammar





Grammar[文法]

- Formal definition[形式化定义]: 4 components {T, N, s, δ}
- T: set of terminal symbols[终结符]
 - Basic symbols from which strings are formed
 - Essentially tokens from lexer leaves in the parse tree
- N: set of non-terminal symbols[非终结符]
 - Each represents a set of strings of terminals internal nodes
 - E.g.: declaration, statement, loop, ...
- s: start symbol[开始符号]
 - One of the non-terminals
- *o*: set of productions[产生式]
 - Specify the manner in which the terminals and non-terminals can be combined to to form strings
 - "LHS \rightarrow RHS": left-hand-side produces right-hand-side



Grammar (cont.)

- Usually, we can just write the σ [简写]
- Merge rules sharing the same LHS[规则合并]
 α → β₁, α → β₂, ..., α → β_n
 α → β₁ | β₂ | ... | β_n







Syntax Analysis[语法分析]

- Informal description of variable declarations in C[变量声明]
 - Starts with int or float as the first token[类型]
 - Followed by one or more *identifier* tokens, separated by token *comma*[逗号分隔的标识符]
 - Followed by token *semicolon*[分号]
- To *check <u>whether a program is well-formed</u>* requires a specification of <u>what is a well-formed program</u>[语法定义]
 - The specification be precise[正确]
 - The specification be complete[完备]
 - Must cover all the syntactic details of the language
 - The specification must be convenient[便捷] to use by both language designer and the implementer
- A context free grammar meets these requirements



