

Programming Assignment II

15-李晓畅-凌之寒

1.1 词法分析器实现与 cool.flex 代码简述

按照给出框架，分为个部分分别实现：

1.运算符和关键字：

按照文档给出的要求：

The keywords of cool are: **class**, **else**, **false**, **fi**, **if**, **in**, **inherits**, **isvoid**, **let**, **loop**, **pool**, **then**, **while**, **case**, **esac**, **new**, **of**, **not**, **true**. Except for the constants **true** and **false**, keywords are case insensitive. To conform to the rules for other objects, the first letter of **true** and **false** must be lowercase; the trailing letters may be upper or lower case.

我们应该对应实现如上所述的关键字和运算符。

运算符只需要对应声明，然后再对应实现方法即可。

```
LE          <=
DARROW      =>
ASSIGN      <-|
121  {DARROW}    { return (DARROW); }
122  {ASSIGN}    { return (ASSIGN); }
123  {LE}        { return (LE); }
```

但值得注意的是对关键字的处理：

一般关键字应该是不区分大小写的，因而在具体实现上可以类似于：

```
(?i:CLASS)    { return (CLASS); }
(?i:ELSE)     { return (ELSE); }
(?i:FI)       { return (FI); }
```

这种样式。

而布尔变量要求首字母必须小写而其他字母可以忽略大小写；识别出对应的布尔后应该将其加入符号表供下一步使用，因而实现上就类似于：

```
TRUE          t[Rr][Uu][Ee]
FALSE         f[Aa][Ll][Ss][Ee]
156  {FALSE} {
157      cool_yylval.boolean = false;
158      return (BOOL_CONST);
159  }
```

2. 整数、标识符和特殊标记

对应定义符合其要求的规则：

```
65  INT          [0-9]+
66  NAMECHAR     [0-9a-zA-Z_]
67  TYPE         [A-Z]{NAMECHAR}*
68  OBJECT       [a-z]{NAMECHAR}*
```

再以类似于下面的方式识别，并加入符号表即可。

```
160  {INT} {
161      cool_yylval.symbol = inttable.add_string(yytext);
162      return (INT_CONST);
163  }
```

3. 注释

单行注释容易实现

```
115  -- .* {}
```

。

实现多行注释需要添加一个状态 `%x COMMENT` 和一个表示注释层数的计数器

```
49  static int commentLayer =0;
```

。

具体来说，正常情况下：在读到`(*`时将开始 `COMMENT` 状态，

```
82  " (*" {
83      BEGIN(COMMENT);
84      commentLayer = 1;
85  }
```

读到大多数字符都可以直接忽略掉；读到(或*)应该修改注释层数，若是最外层的*)就应该结束注释状态；

```
105  <COMMENT>"*)" {
106      commentLayer --;
107      if (commentLayer == 0 ) BEGIN(0);
108  }
```

读到\n 应该增加行数。

```
87  <COMMENT>"\n" {
88      curr_lineno ++;
89  }
```

特殊地，注释中包含 EOF 符将报告 EOF in comment 错；未启动 COMMENT 状态下读到*)应该报 Unmatched *)错。

4.字符串

使用一个字符串作为缓冲记录识别到的字符

```
50  static std::string str="";
```

用 STR 状态和 ERRORSTR 状态标记字符串正在识别字符串和发现错误字符串

```
73  %X STR
74  %X ERRORSTR。
```

具体来说，正常情况下识别到前引号” 将开始识别字符串过程；

字符串识别过程对大多数字符而言都是直接加入到缓冲，但对表达转义的字符应该特殊处理；

```
199  <STR>\\. {
200      // std :: cout<<"SPACE: "<<yytext<<endl;
201      switch (yytext[1]){
202          case 'n': str.append("\n"); break;
203          case 'b': str.append("\b"); break;
204          case 'f': str.append("\f"); break;
205          case 't': str.append("\t"); break;
206          default : str.push_back(yytext[1]);
207      }
208  }
```

可以换行表示字符串，但必须加\；

```
219  <STR>\\\n {
220      str.push_back(yytext[1]);
221      curr_lineno ++;
222  }
```

在识别到后引号且长度符合定义即可结束对该字符串的识别。

特殊地，单个的\0 在字符串中总是非法的，将报包含空字符错；

```
225  <STR>. {
226      if (yytext[0]==0) {cool_yylval.error_msg = "String contains null character."; BEGIN(ERRORSTR);}
227      str.push_back(yytext[0]);
228  }
```

对于上面描述的错误，将进入 `ERRORSTR` 状态，此状态下什么都不做直到找到可以重启的地方再次开始识别；

```
237 <ERRORSTR>["\n] {
238 | BEGIN(0);
239 | return ERROR;
240 }
```

COOL 语言不允许字符串直接换行：

```
229 <STR>[\n] {
230 | cool_yylval.error_msg = "Unterminated string constant";
231 | curr_lineno ++;
232 | BEGIN(0);
233 | return ERROR;
234 }
```

字符串未正常结束情况下读到 EOF 将报对应报错；

```
193 <STR><<EOF>> {
194 | cool_yylval.error_msg = "EOF in string constant";
195 | BEGIN(0);
196 | return ERROR;
197 }
```

即使字符串正常识别到后引号，仍有可能因过长而出错：

```
210 <STR>\" {
211 | if (str.length() >= MAX_STR_CONST ) {
212 | | cool_yylval.error_msg="String constant too long";
213 | | BEGIN(0); return ERROR;}
214 | BEGIN(0);
215 | cool_yylval.symbol = stringtable.add_string((char*)str.c_str());
216 | return (STR_CONST);
217 }
```

5.其他字符

一般状态下，对读到的空格符直接忽略就可以了

```
63 SPACE [ \t\f\r\v] 175 {SPACE} {}
```

COOL 支持的合法字符集并不很多，可以将其单独写出来：

```
245 [\\(\\)\\{\\}<=,,:;~@] {
246 | | return yytext[0];
247 | }
248
249 [\\-\\+\\*\\/\\. ] {
250 | | return yytext[0];
251 | }
```

那对于剩下的其他字符应该对应报错：

```
253 . {
254 | if (yytext[0]==0) {cool_yylval.error_msg = "\\000"; return ERROR;}
255 | cool_yylval.error_msg = yytext;
256 | return ERROR;
257 | }
```

1.2 验证

按照要求执行代码，结果如下：

```
● wojtek@wojtek-virtual-machine:~/user/Compile-Lab/assignments/PA2$ make lexer
flex -d -ocool-lex.cc cool.flex
/bin/sh -ec 'g++ -MM -I. -I../include/PA2 -I../src/PA2 cool-lex.cc | sed '\''s/(cool-lex.o)[ :]*\1 cool-lex.d : /g'\'' > cool-lex.d'
/bin/sh -ec 'g++ -MM -I. -I../include/PA2 -I../src/PA2 handle_flags.cc | sed '\''s/(handle_flags.o)[ :]*\1 handle_flags.d : /g'\'' > handle_flags.d'
/bin/sh -ec 'g++ -MM -I. -I../include/PA2 -I../src/PA2 stringtab.cc | sed '\''s/(stringtab.o)[ :]*\1 stringtab.d : /g'\'' > stringtab.d'
/bin/sh -ec 'g++ -MM -I. -I../include/PA2 -I../src/PA2 utilities.cc | sed '\''s/(utilities.o)[ :]*\1 utilities.d : /g'\'' > utilities.d'
/bin/sh -ec 'g++ -MM -I. -I../include/PA2 -I../src/PA2 lextest.cc | sed '\''s/(lextest.o)[ :]*\1 lextest.d : /g'\'' > lextest.d'
g++ -g -Wall -Wno-unused -Wno-write-strings -I. -I../include/PA2 -I../src/PA2 -c cool-lex.cc
g++ -g -Wall -Wno-unused -Wno-write-strings -I. -I../include/PA2 -I../src/PA2 lextest.o utilities.o stringtab.o handle_flags.o cool-lex.o -lfl -o lexer
● wojtek@wojtek-virtual-machine:~/user/Compile-Lab/assignments/PA2$ ./checker test.cl
passed
```

同预期一致，测试通过。

2.测试样例构建

对 4.1 和 4.4 要求的错误构造样例：

```
1  `~!@#$$%^&*()_+=[ ]'<>?/.
2  __what_am_I
3  LET_STMT --This is strange! *)
4  "I am an unescaped
5  newline"
6  "I am an escaped\
7  newline"
8  *)
9  (*I am EOF in comment
```

包含对合法/非法字符；变量名；单独出现的*)；unescaped/escaped 换行和 EOF 的检查。

同预期一致

```
● wojtek@wojtek-virtual-machine:~/user/Compile-Lab/assignments/PA2$ ./checker test2
passed
```

对 4.3 构造样例：

```
1  "123\0"
2  "123\\0"
3  "123\\\0"
4  "123\\\\0"
5  "123164564156156165149874896165498165847986151231645641561561651498748961654981658479861512316456415615"
6  "123\nab"
7  "
```

包含对字符串中转义字符、超长字符串和 EOF 的检查。

同预期一致

```
● wojtek@wojtek-virtual-machine:~/user/Compile-Lab/assignments/PA2$ ./checker test1
passed
```

3.特殊情况和处理

我们发现，如果在正文中出现\0 标准词法分析器会报告 #6 ERROR "\000"。

我们将其作为特殊情况处理。 `if (yytext[0]==0) {cool_yylval.error_msg = "\000"; return ERROR;}`

4.实验总结

本次实验中我们运用 flex 工具实现了 COOL 语言的词法分析器。

flex 将符合指定格式的文件, 转化为 C 代码, 用于进行**词法分析** Lexical Analysis。在.flex 文件中, 我们通过设置**正则表达式**, 定义了一些**词汇** lexeme, 这些词汇将被转化为抽象的符号 token, 作为词法分析的结果。

在词法分析器的设计中我们不难看出 COOL 语言一些独特安排的好处。比如对 type 和 object 的特殊要求、提供较少的运算符和其他符号等设计可以使得我们在构建词法分析器时更简单方便。

良好的词法分析器需要精巧的设计, 正确地规划先后顺序、长度才能正确实现功能并保证正确。

附录

Cool.flex 代码如下:

```
1.  /*
2.   *   The scanner definition for COOL.
3.   */
4.
5.  /*
6.   *   Stuff enclosed in %{ %} in the first section is copied verbatim to the
7.   *   output, so headers and global definitions are placed here to be visible
8.   *   to the code in the file.  Don't remove anything that was here initially
9.   */
10. %{
11. #include <cool-parse.h>
12. #include <stringtab.h>
13. #include <utilities.h>
14. #include <stdint.h>
15. #include <string>
16. #include <iostream>
17.
18. /* The compiler assumes these identifiers. */
19. #define yylval cool_yylval
20. #define yylex   cool_yylex
21.
22. /* Max size of string constants */
23. #define MAX_STR_CONST 1025
24. #define YY_NO_UNPUT    /* keep g++ happy */
25.
26. extern FILE *fin; /* we read from this file */
27.
28. /* define YY_INPUT so we read from the FILE fin:
29.  * This change makes it possible to use this scanner in
30.  * the Cool compiler.
31.  */
32. #undef YY_INPUT
33. #define YY_INPUT(buf,result,max_size) \
34.     if ( (result = fread( (char*)buf, sizeof(char), max_size, fin)) < 0) \
35.         YY_FATAL_ERROR( "read() in flex scanner failed");
36.
37. char string_buf[MAX_STR_CONST]; /* to assemble string constants */
38. char *string_buf_ptr;
39.
40. extern int curr_lineno;
41. extern int verbose_flag;
42.
43. extern YYSTYPE cool_yylval;
44.
45. /*
```

```
46.  *   Add Your own definitions here
47.  */
48.
49. static int commentLayer =0;
50. static std::string str="";
51.
52. static int errorflag = 0;
53.
54. %}
55.
56. /*
57.  * Define names for regular expressions here.
58.  */
59. LE          <=
60. DARROW      =>
61. ASSIGN      <-
62.
63. SPACE       [ \t\f\r\v]
64.
65. INT         [0-9]+
66. NAMECHAR    [0-9a-zA-Z_]
67. TYPE        [A-Z]{NAMECHAR}*
68. OBJECT      [a-z]{NAMECHAR}*
69. TRUE        t[Rr][Uu][Ee]
70. FALSE       f[Aa][Ll][Ss][Ee]
71.
72. %x COMMENT
73. %x STR
74. %x ERRORSTR
75.
76. %option noyywrap
77. %%
78. /*
79.  *   Nested comments
80.  */
81.
82. "(" {
83.     BEGIN(COMMENT);
84.     commentLayer = 1;
85. }
86.
87. <COMMENT>"\n" {
88.     curr_lineno ++;
89. }
90.
91. <COMMENT><<EOF>> {
92.     cool_yylval.error_msg = "EOF in comment";
93.     BEGIN(0);
94.     return ERROR;
95. }
96.
97. <COMMENT>. {
98.
99. }
100.
101. <COMMENT>"(" {
102.     commentLayer ++;
103. }
104.
105. <COMMENT>*)" {
106.     commentLayer --;
107.     if (commentLayer == 0 ) BEGIN(0);
108. }
109.
110. \*\) {
111.     cool_yylval.error_msg = "Unmatched *");
112.     return (ERROR);
113. }
114.
```

```
115.  --.* {}
116.
117.
118.  /*
119.   * The multiple-character operators.
120.   */
121.  {DARROW}      { return (DARROW); }
122.  {ASSIGN}      { return (ASSIGN); }
123.  {LE}          { return (LE); }
124.
125.  /*
126.   * Keywords are case-insensitive except for the values true and false,
127.   * which must begin with a lower-case letter.
128.   */
129.  [\[\]\\'>] {
130.      cool_yylval.error_msg = yytext;
131.      return (ERROR);
132.  }
133.  (?i:CLASS)     { return (CLASS); }
134.  (?i:ELSE)      { return (ELSE); }
135.  (?i:FI)        { return (FI); }
136.  (?i:IF)        { return (IF); }
137.  (?i:IN)        { return (IN); }
138.  (?i:INHERITS)  { return (INHERITS); }
139.  (?i:LET)       { return (LET); }
140.  (?i:LOOP)      { return (LOOP); }
141.  (?i:POOL)      { return (POOL); }
142.  (?i:THEN)      { return (THEN); }
143.  (?i:WHILE)     { return (WHILE); }
144.  (?i:CASE)      { return (CASE); }
145.  (?i:ESAC)      { return (ESAC); }
146.  (?i:OF)        { return (OF); }
147.  (?i:NEW)       { return (NEW); }
148.  (?i:LE)        { return (LE); }
149.  (?i:NOT)       { return (NOT); }
150.  (?i:ISVOID)    { return (ISVOID); }
151.
152.  {TRUE} {
153.      cool_yylval.boolean = true;
154.      return (BOOL_CONST);
155.  }
156.  {FALSE} {
157.      cool_yylval.boolean = false;
158.      return (BOOL_CONST);
159.  }
160.  {INT} {
161.      cool_yylval.symbol = inttable.add_string(yytext);
162.      return (INT_CONST);
163.  }
164.
165.  {TYPE} {
166.      cool_yylval.symbol = idtable.add_string(yytext);
167.      return (TYPEID);
168.  }
169.
170.  {OBJECT} {
171.      cool_yylval.symbol = idtable.add_string(yytext);
172.      return (OBJECTID);
173.  }
174.
175.  {SPACE} {}
176.
177.  [\n] {
178.      curr_lineno ++;
179.  }
180.
181.  /*
182.   * String constants (C syntax)
183.   * Escape sequence \c is accepted for all characters c. Except for
```

```

184.     * \n \t \b \f, the result is c.
185.     *
186.     */
187.
188.     \" {
189.         BEGIN (STR);
190.         str = "";
191.     }
192.
193.     <STR><<EOF>> {
194.         cool_yylval.error_msg = "EOF in string constant";
195.         BEGIN(0);
196.         return ERROR;
197.     }
198.
199.     <STR>\\. {
200.         // std :: cout<<"SPACE: "<<yytext<<endl;
201.         switch (yytext[1]){
202.             case 'n': str.append("\n"); break;
203.             case 'b': str.append("\b"); break;
204.             case 'f': str.append("\f"); break;
205.             case 't': str.append("\t"); break;
206.             case '\\0': cool_yylval.error_msg = "String contains escaped null character."; BEGIN(ERRORSTR); break;
207.             default : str.push_back(yytext[1]);
208.         }
209.     }
210.
211.     <STR>\" {
212.         // cout<<"DONE "<<str.length();
213.         if (str.length()>= MAX_STR_CONST ) {cool_yylval.error_msg="String constant too long"; BEGIN(0); return ERROR;}
214.         // cout<<"New str: "<<str<<endl;
215.         BEGIN(0);
216.         cool_yylval.symbol = stringtable.add_string((char*)str.c_str());
217.         return (STR_CONST);
218.     }
219.
220.     <STR>\\\n {
221.         str.push_back(yytext[1]);
222.         curr_lineno ++;
223.     }
224.
225.     <STR>. {
226.         if (yytext[0]==0) {cool_yylval.error_msg = "String contains null character."; BEGIN(ERRORSTR);}
227.         str.push_back(yytext[0]);
228.     }
229.
230.     <STR>[\n] {
231.         cool_yylval.error_msg = "Unterminated string constant";
232.         curr_lineno ++;
233.         // cout<< "EOL in string constant"<< endl;
234.         BEGIN(0);
235.         return ERROR;
236.     }
237.
238.     <ERRORSTR>[\"\\n] {
239.         BEGIN(0);
240.         return ERROR;
241.     }
242.
243.     <ERRORSTR>. {
244.
245.     }
246.
247.     [\\(\\)\\{\\}<=,,:;~@] {
248.         return yytext[0];
249.     }
250.
251.     [\\-\\+\\*\\/\\.] {
252.         return yytext[0];

```

```
253.  }
254.
255.  . {
256.      if (yytext[0]==0) {cool_yylval.error_msg = "\000"; return ERROR;}
257.      cool_yylval.error_msg = yytext;
258.      return ERROR;
259.  }
260.
261.  %%
```