# 实验三、词法分析实验

学号：*1120180207*

姓名：*唐小娟*

班级：*07111801*

## 1. 实验目的

（1）熟悉 C 语言的词法规则，了解编译器词法分析器的主要功能和实现技术，掌握典型词法分析器构造方法，设计并实现 C 语言词法分析器；

（2）了解 Flex 工作原理和基本思想，学习使用工具自动生成词法分析器；

（3）掌握编译器从前端到后端各个模块的工作原理，词法分析模块与其他模块之间的交互过程。

## 2. 实验内容

根据 C 语言的词法规则，设计识别 C 语言所有单词类的词法分析器的确定有限状态自动机，并使用 Java、C\C++或者 Python 其中任何一种语言，采用程序中心法或者数据中心法设计并实现词法分析器。

词法分析器的输入为 C 语言源程 序，输出为属性字流。

学生可以选择编码实现词法分析器，也可以选择使用 Flex 自动生成词法分析器。需要注意的是，Flex 生成的是 C 为实现语言的词法分析器，如果需要生成 Java 为实现语言的词法分析器，可以尝试 JFlex 或者 ANTLR。 由于框架是基于 Java 语言实现的，并且提供了相应的示例程序，建议学生 使用 Java 语言在示例的基础上完成词法分析器。

本实验我采用Flex来生成以C为实现语言的词法生成器，其生成的代码文件为lex.yy.c，之后用gcc进行编译得到词法生生成器的可执行文件myscanner.exe。

## 3. 实验环境

| 名称 | 信息 |
|---|---|
| 操作系统版本 | Ubuntu 20.04.1 LTS、Windows家庭中文版 |
| Flex版本 | 2.6.4 |
| GCC版本 | 8.1.0 |
| gedit版本 | 3.36.2 |

## 4. 实验过程

　　该实验以 C 语言作为源语言，构建 C 语言的词法分析器，对于给定的测试程 序，输出属性字符流。词法分析器的构建按照 C 语言的词法规则进行。C 语言的 发展经历了不同的阶段，早期按照 C99 标准进行编程和编译器的实现，2011 年 又对 C 语言规范进行了修订，形成了 C11（又称 C1X）。下面以 C11 为基准，对 C 语言的词法规则进行简要的描述和相关设计。

## 4.1 设计思路

- 对语言的各类单词分别构造状态图；
- 将各类状态图进行合并，构成一个能识别该语言所有单词的状态图；
- 将各类单词的状态图的初态合并为唯一初态；
- 调整冲突的编号

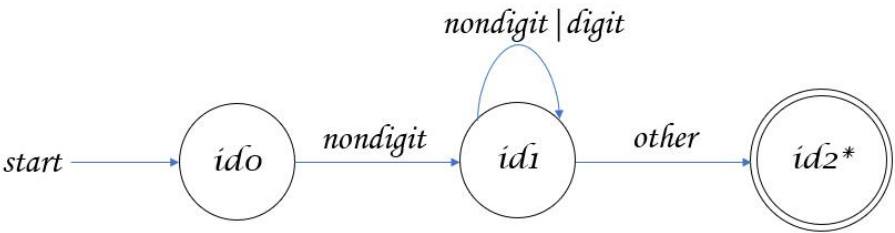　　在这里，我分成了三类单词：标识符、整数常量、浮点数常量、字符常量、字符串常量、运算符和界限符。为了简化实现过程，具体的定义以老师给定文档为标准。

### 4.1.1标识符

　　定义：只能由字母、数字和下划线三种字符组成，且第一个字符必须是字母或下划线。

C 语言标识符的定义如下：

```
identifier →     identifier-nondigit
                | identifier identifier-nondigit
                | identifier digit
identifier-nondigit →    nodigit
                       | universal-character-name
                       | other implementation-define characters

nondigit →  _ | a | b | c | d | e | f | g | h | i | j | k | l | m | n | o | p | q | r
            | s | t | u | v | w | x | y | z | A | B | C | D | E | F | G | H | I
            | J | K | L | M | N | O | P | Q | R | S | T | U | V | W | X | Y
            | Z
digit    →  0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
```

DFA：



　　其中，other代表非nondigit|digit的字符，状态2为结束状态，且代表回退一个字符。除此之外，由于关键字也满足上述DFA，所以还需要进行关键字的筛选。

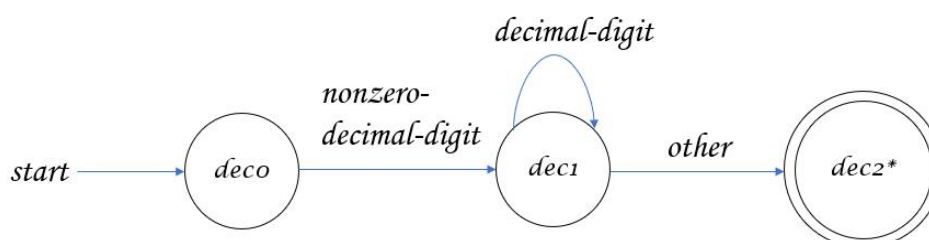| | | | | |
|---|---|---|---|---|
| *auto* | *break* | *case* | *char* | *const* |
| *continue* | *default* | *do* | *double* | *else* |
| *enum* | *extern* | *float* | *for* | *goto* |
| *if* | *inline* | *int* | *long* | *register* |
| *restrict* | *return* | *short* | *signed* | *sizeof* |
| *static* | *struct* | *switch* | *typedef* | *union* |
| *unsigned* | *void* | *volatile* | *while* | |

### 4.1.2 整形常量

定义：C语言的整形常量主要包括：十进制、八进制、十六进制。依照这三类的定义分别实现DFA。
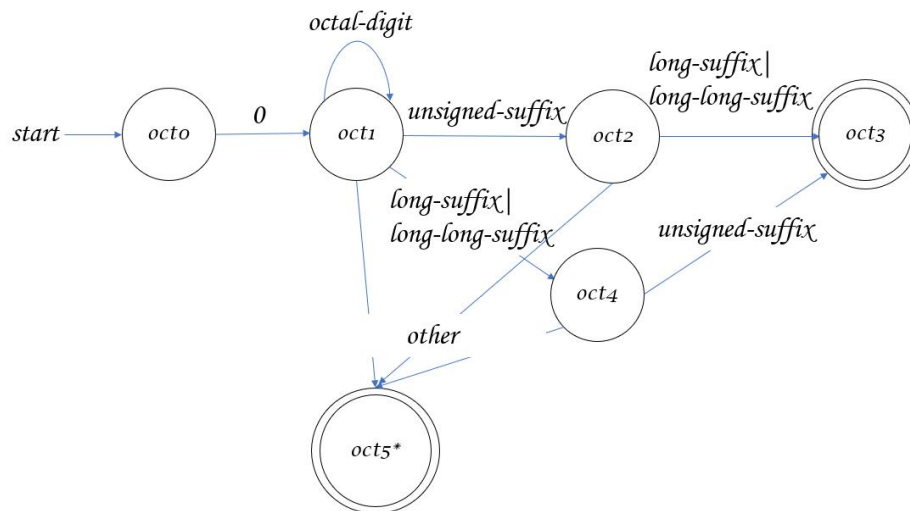
C语言整型常量的定义如下：

*integer-constant* → *decimal-constant integer-suffix*
　　　　　　　　| *octal-constant integer-suffix*
　　　　　　　　| *hexadecimal-constant integer-suffix*
*decimal-constant* → *nonzero-digit* | *decimal-constant digit*
*octal-constant* → *0* | *octal-constant octal-digit*
*hexadecimal-constant* → *hexadecimal-prefix hexadecimal-digit*
　　　　　　　　| *hexadecimal-constant hexadecimal-digit*
*hexadecimal-prefix* → *0x* | *0X*
*nonzero-digit* → *1* | *2* | *3* | *4* | *5* | *6* | *7* | *8* | *9*
*octal-digit* → *0* | *1* | *2* | *3* | *4* | *5* | *6* | *7*
*hexadecimal-digit* → *0* | *1* | *2* | *3* | *4* | *5* | *6* | *7* | *8* | *9* | *a* | *b* | *c* | *d* | *e* | *f* |
　　　　　　　　*A* | *B* | *C* | *D* | *E* | *F*
*integer-suffix* → *unsigned-suffix long-suffix*
　　　　　　　　| *unsigned suffix long-long suffix*
　　　　　　　　| *long-suffix unsigned-suffix*
　　　　　　　　| *long-long-suffix unsigned-suffix*
*unsigned-suffix* → *u* | *U*
*long-suffix* → *l* | *L*
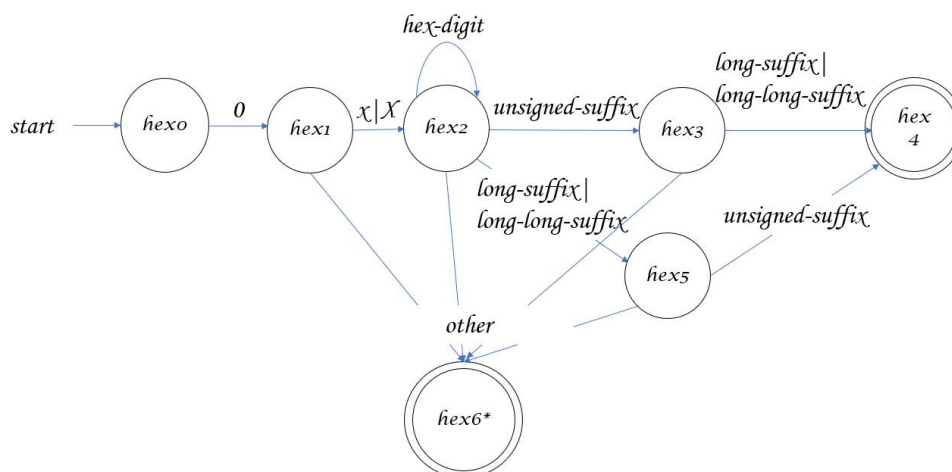*long-long-suffix* → *ll* | *LL*

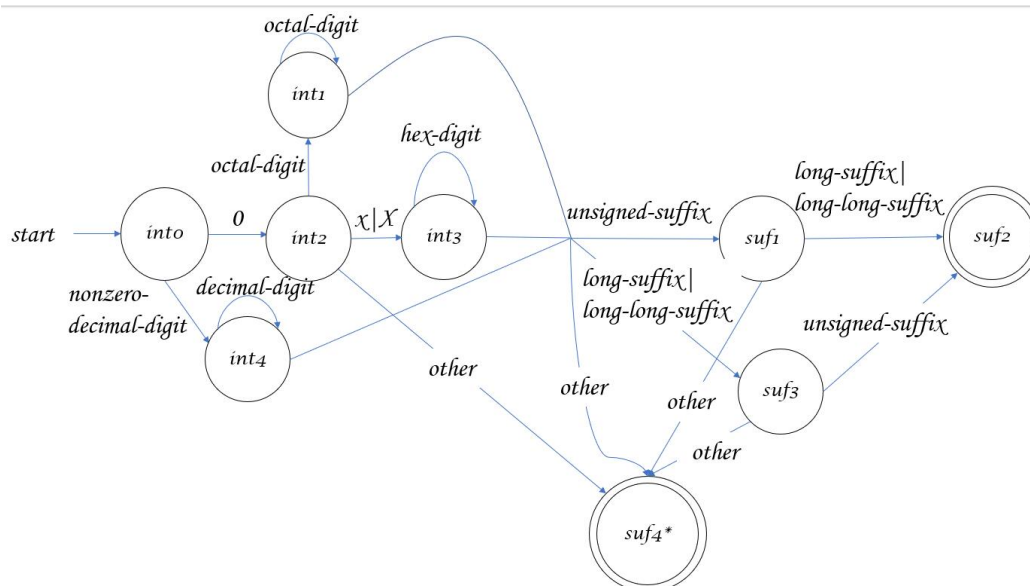**decimal-constant**

## octal-constant



## hexadecimal-constant



所以，将这三个DFA结合起来，有Interger的DFA：
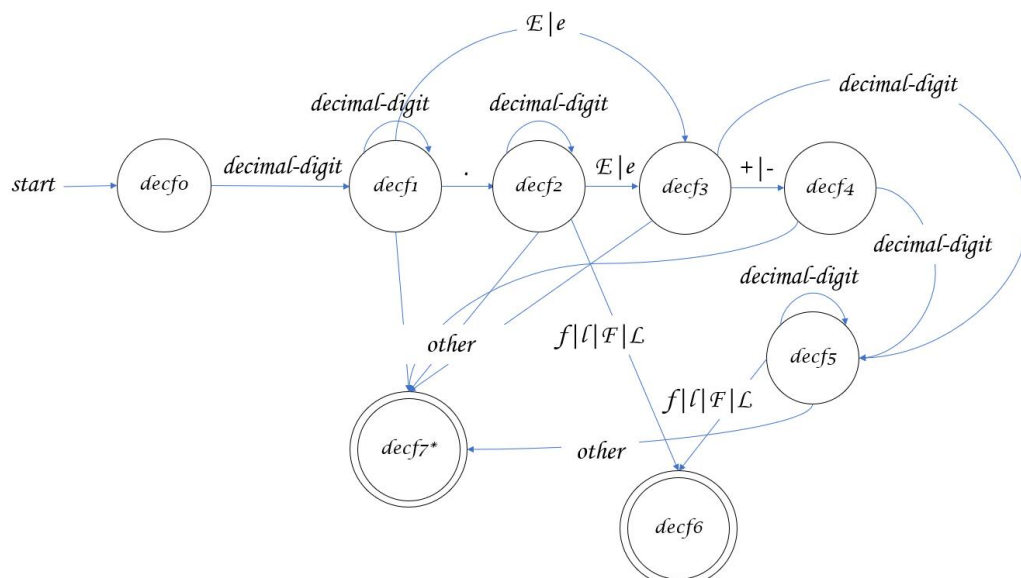
### 4.1.3 浮点数常量

定义：C语言的浮点数常量主要包括：十进制、十六进制。依照这三类的定义分别实现DFA。
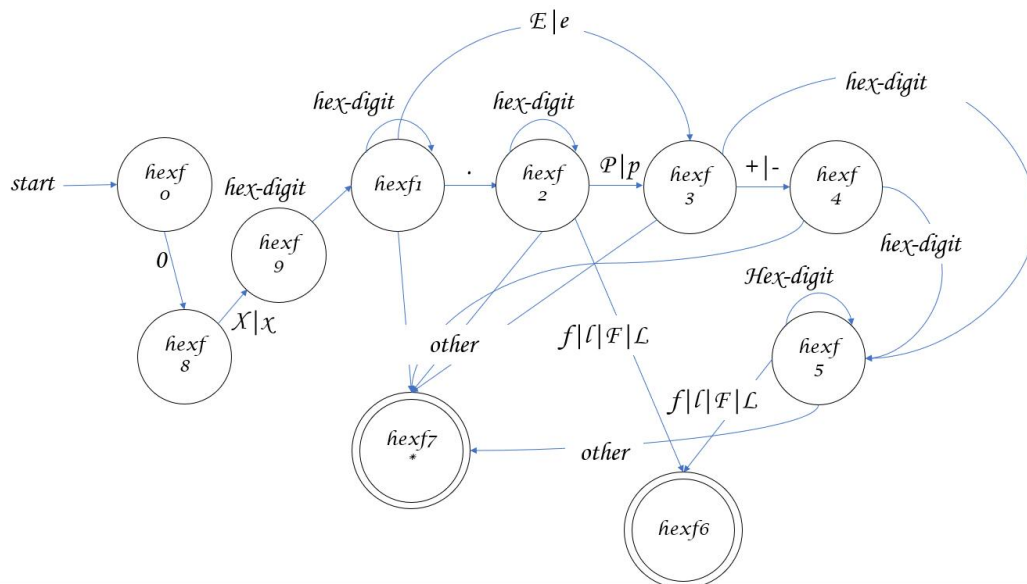
C 语言浮点型常量定义如下：

```
floating-constant→ decimal-floating-constant
                 | hexadecimal-floating-constant
decimal-floating-constant    →   fractional-constant exponent-part floating-suffix
                                | digit-sequence exponent-part floating-suffix
hexadecimal-floating-constant    →   hexadecimal-prefix hexadecimal-fraction-constant
                                        binary-exponent-part floating-suffix
                                    | hexadecimal-prefix hexadecimal-digit-sequence
                                        binary-exponent-part floating-suffix
fractional-constant →   digit-sequence . digit-sequence | digit-sequence .
exponent-part      →   e sign digit-sequence | E sign digit-sequence
sign               →   + | -
digit-sequence     → digit | digit-sequence digit
hexadecimal-fractional-constant → hexadecimal-digit-sequence . hexadecimal-digit-sequence
                        | hexadecimal-digit-sequence .
binary-exponent-part → p sign digit-sequence | P sigh digit-sequence
hexadecimal-digit-sequence → hexadecimal-digit
                        | hexadecimal-digit-sequence hexadecimal-digit
floating-suffix → f | l | F | L
```

**decimal-floating-constant**



**hexadecimal-floating-constant**

所以将这两个DFA结合起来，得到float的DFA：



除 f0、f13* 外，其他状态均有一条other 线连入

### 4.1.4 字符常量

经过查阅资料，可知字符可以有多个字符，如'abcd'，虽然会有warning，但编译执行仍然是通过的。在这里为了简便，就不考虑出错处理了。简单的认为单引号引起的内容就是字符常量。

C 语言字符常量定义如下：

*character-constant→ ' c-char-sequence ' | L' c-char-sequence ' | u' c-char-sequence '*
             *| U' c-char-sequence '*
*c-char-sequence → c-char | c-char-sequence c-char*
*c-char → any member of the source character set except the single-quote ', backslash \, or new-line character*
        *| escape-sequence*
*escape-sequence → simple-escape-sequence | octal-escape-sequence*
             *| hexadecimal-escape-sequence | universal-character-name*
*simple-escape-sequence → \' | \" | \? | \\ | \a | \b | \f | \n | \r | \t | \v*

DFA为：



### 4.1.5 字符串常量

定义如下：

C 语言字符串字面量定义如下：

*string-literal→encoding-prefix " s-char-sequence "*
*encoding-prefix → u8 | u | U | L*
*s-char-sequence → s-char | s-char-sequence s-char*
*s-char → any member of the source character set except the double-quote ", backslash \, or new-line character*
        *| escape-sequence*

DFA：

### 4.1.6 运算符和界限符

根据文档定义的运算符和界限符内容：由于<%、%>、%:、%:%:、...并不常用，所以就忽略了。

| [ | ] | ( | ) | { | } | . | -> |
|---|---|---|---|---|---|---|---|
| ++ | -- | & | * | + | - | ~ | ! |
| / | % | << | >> | < | > | <= | >= |
| == | != | ^ | \| | && | \|\| | | |
| ? | : | ; | ... | | | | |
| = | *= | /= | %= | += | -= | <<= | |
| >>= | | | | | | | |
| &= | ^= | \|= | , | # | ## | <: | :> |
| <% | %> | %: | %:%: | | | | |

DFA如下：

除 sig0、sig1、sig3、sig9、sig13 结束状态外，其他所有状态均有一条 other 线连到 (sig15) 状态。

### 4.1.7 总DFA

将前文所有设计的 DFA 合并到一起称为最终的 DFA，由于汇总起来得到的DFA实在过于庞大，因此将这个整个大的DFA先分成五个子图。如下：（具体的不再示意）

DFA1包括了标识符和关键词的部分；
DFA2包括了浮点数和整形数
DFA3包括了运算符和界限符
DFA4包括了字符串常量
DFA5包括了字符常量

## 4.2 代码实现

本次实验采取Flex自动生成词法分析器，关键部分在于上述单词正则表达式的书写，除此之外，在匹配的过程中，从上到下优先级从高到低，也就是说，先匹配到上方的正则表达式之后，就不会再进行匹配了。例如：匹配关键词要先于标识符。myscanner.l具体部分关键代码如下：

```
1   KEYWORD (auto)|(break)|(case)|(char)|(const)|(continue)|
        (default)|(do)|(double)|(else)|(enum)|(extern)|(float)|(for)|
        (goto)|(if)|(inline)|(int)|(long)|(register)|(restrict)|
        (return)|(short)|(signed)|(sizeof)|(static)|(struct)|(switch)|
        (typedef)|(union)|(unsigned)|(void)|(volatile)|(while)
2
3   IDENTIFIER ([_a-zA-Z][_a-zA-Z0-9]*)
4
5   THREESYMBOL [<][<][=]|[>][>][=]
6
7   DOUBLESIMBOL [\*\+\^\|=<>!&%-][=]|[-][>]|[\+][\+]|[-][-]|[<]
        [<]|[>][>]|[&][&]|[\|][\|]|[#][#]
8
9   SIMBOL  ([\[\]](),.;:!#{}\-?^<>&~|%*+/=])
10
11  CHAR    ([cLuU]?)(\'.*\')
12
13  STRING  (([uUL]?)|(u8))(\"(.*)\")
14
15  INTERGER (([1-9][0-9]*)|(0[0-7]*)|(0[xX][0-9a-fA-F]*))([uU]?
        (ll|LL)?|(ll|LL)?[uU]?|[uU]?[Ll]?|[Ll]?[Uu]?)
16
17  DECIFLOAT ([0-9]+\.?[0-9]*)([eE][+-]?[0-9]+)?([flFL])?
```

```
18   HEXAFLOAT  (0[xX][0-9a-fA-F]+\.?[0-9a-fA-F]*)([pP][+-]?[0-9]+)?
     ([flFL])?
19   COMMENT    (\/\/.*)
20
21   {KEYWORD}        {
22
23                       numStart = temp+1;
24                       numEnd=numStart + strlen(yytext)-1;
25                       temp = numEnd;
26                       numTokens++;
27                       return T_keyword;
28                   }
29
30   //...
31
32
33   int main(int argc,char* argv[])
34   {
35       int token_type;
36
37       yyin=fopen(argv[1],"r");
38       yyout=fopen(argv[2],"w");
39       printf("%s\n",argv[1]);
40       printf("%s\n",argv[2]);
41       while (token_type = yylex()) {
42           if(token_type == T_empty) continue;
43           if(token_type == T_symbol)
44               fprintf(yyout,"[@%d,%d:%d=\'%s\',
     <'%s'>,%d:%d]\n",numTokens,numStart,numEnd,yytext,yytext,numLin
     es,numStart);
45           else
46               fprintf(yyout,"[@%d,%d:%d=\'%s\',
     <%s>,%d:%d]\n",numTokens,numStart,numEnd,yytext,token_strs[toke
     n_type],numLines,numStart);
47       }
48       return 0;
49   }
```

myscanner.h的接口文件（方便输出相关内容）：

```
1   #ifndef MYSCANNER_H
2   #define MYSCANNER_H
3
4
5
6
7   typedef enum{
8
    T_over,T_keyword,T_identifier,T_interger,T_float,T_char,T_strin
    g,T_symbol,T_empty
9
```

```
10   }TokenType;
11
12   const char* token_strs[]={
13
     "EOF","keyword","Identifier","const_interger","const_float","co
     nst_char","const_string"
14   };
15
16
17   #endif
18
```

之后运行flex myscanner.l生成的lex.yy.c源文件，再通过gcc编译生成可执行文件，修改老师给定的BIT-MiniCC框架中的config.xml文件（type和path），运行之后，可以得到相应的tokens文件。

```
Start to compile ...
D:\Programming\eclipse\workplace\bitmincc-clean\test\scan_test\1_scanner_test.c
D:\Programming\eclipse\workplace\bitmincc-clean\test\scan_test\1_scanner_test.tokens
Compiling completed!
```

## 5. 实验结果

根据老师给定的测试文件1_scanner_test.c文件

```
1    int main()
2    {
3        //integer-constant
4        +1 ;
5        -10 ;
6        1000 ;
7        1l ;
8        1U ;
9        10ul ;
10       10LU ;
11       1000ull ;
12       1000LLU ;
13
14       +0 ;
15       -00 ;
16       007 ;
17       00ul ;
18       00LLU ;
19
20       +0x0 ;
21       -0x00 ;
22       0XABCDEF ;
23       0xful ;
24       0XFLLU ;
25
```

```
26    //floating-constant
27        0.0 ;
28        +1.1e+1 ;
29        -1.1E-1 ;
30        1.1e1f ;
31        1.1E1L ;
32
33        0x0p0;
34        0x0.0p0 ;
35        +0xa.ap+1 ;
36        -0xa.aP-1 ;
37        0xa.ap1f ;
38        0xa.aP1L ;
39
40    //character-constant
41        'a' ;
42        L'a' ;
43        U'a' ;
44        u'a' ;
45        '\n' ;
46        '\?' ;
47        '\24' ;
48
49    //string-literal
50        "abcdefg123456\\" ;
51        u8"a" ;
52        u"a" ;
53        U"a" ;
54        L"a" ;
55
56    //identifier
57        int __a ;
58        int __1 ;
59        int a1 ;
60        int a ;
61        int a_1 ;
62        int a_ ;
63
64    //keyword
65        int i=1 ;
66        float f ;
67        double d ;
68        char c ;
69        long l ;
70        short s ;
71        signed si ;
72        unsigned short us ;
73        typedef struct
74        {
75            int num1;
76        }test;
77        test test1 ;
```

```c
78        test *test2 ;
79        static int sti ;
80        const int ci ;
81        sizeof( int ) ;
82        if( 1 )
83        {

85        }
86        else
87        {

89        }
90        for( ; ; )
91        {
92            continue ;
93        }
94        while( 1 )
95        {

97        }
98        switch(i){
99            case 1:  break;
100            default:  break;
101        }
102        do
103        {

105        }while( 0 ) ;
106        goto gotoFlag ;
107        //operators
108        int array[10] = {0} ;
109        test1.num1 = 0 ;
110        test2->num1 = 0 ;
111        i++ ;
112        i-- ;
113        i = 1 + 1 ;
114        i = 1 - 1 ;
115        i = 1 * 1 ;
116        i = 1 / 1 ;
117        i = 1 % 1 ;
118        i = !i ;
119        i = i & 1 ;
120        i = i | 1 ;
121        i = i && 1 ;
122        i = i || 1 ;
123        i = i ^ 1 ;
124        i = i >> 1 ;
125        i = i << 1 ;
126        i = i > 1 ? 1 : 2 ;
127        i += 1 ;
128        i -= 1 ;
129        i *= 1 ;
```

```
130        i /= 1 ;
131        i %= 1 ;
132        i &= 1 ;
133        i |= 1 ;
134        i ^= 1 ;
135        i >>= 1 ;
136        i <<= 1 ;
137        if( i==1 ){}
138        if( i!=1 ){}
139        if( i>1 ){}
140        if( i<1 ){}
141        if( i>=1 ){}
142        if( i<=1 ){}
143    gotoFlag:
144        return 0 ;
145    }
146    void function( int arg1 , int arg2 )
147    {
148    }
```

得到我们的词法分析结果：

```
 1    [@0,0:2='int',<keyword>,1:0]
 2    [@1,4:7='main',<Identifier>,1:4]
 3    [@2,8:8='(',<'('>,1:8]
 4    [@3,9:9=')',<')'>,1:9]
 5    [@4,0:0='{',<'{'>,2:0]
 6    [@5,1:1='+',<'+'>,4:1]
 7    [@6,2:2='1',<const_interger>,4:2]
 8    [@7,4:4=';',<';'>,4:4]
 9    [@8,1:1='-',<'-'>,5:1]
10    [@9,2:3='10',<const_interger>,5:2]
11    [@10,5:5=';',<';'>,5:5]
12    [@11,1:4='1000',<const_interger>,6:1]
13    [@12,6:6=';',<';'>,6:6]
14    [@13,1:2='1l',<const_interger>,7:1]
15    [@14,4:4=';',<';'>,7:4]
16    [@15,1:2='1U',<const_interger>,8:1]
17    [@16,4:4=';',<';'>,8:4]
18    [@17,1:4='10ul',<const_interger>,9:1]
19    [@18,6:6=';',<';'>,9:6]
20    [@19,1:4='10LU',<const_interger>,10:1]
21    [@20,6:6=';',<';'>,10:6]
22    [@21,1:7='1000ull',<const_interger>,11:1]
23    [@22,9:9=';',<';'>,11:9]
24    [@23,1:7='1000LLU',<const_interger>,12:1]
25    [@24,9:9=';',<';'>,12:9]
26    [@25,1:1='+',<'+'>,14:1]
27    [@26,2:2='0',<const_interger>,14:2]
28    [@27,4:4=';',<';'>,14:4]
29    [@28,1:1='-',<'-'>,15:1]
```

```
30    [@29,2:3='00',<const_interger>,15:2]
31    [@30,5:5=';',<';'>,15:5]
32    [@31,1:3='007',<const_interger>,16:1]
33    [@32,5:5=';',<';'>,16:5]
34    [@33,1:4='00ul',<const_interger>,17:1]
35    [@34,6:6=';',<';'>,17:6]
36    [@35,1:5='00LLU',<const_interger>,18:1]
37    [@36,7:7=';',<';'>,18:7]
38    [@37,1:1='+',<'+'>,20:1]
39    [@38,2:4='0x0',<const_interger>,20:2]
40    [@39,6:6=';',<';'>,20:6]
41    [@40,1:1='-',<'-'>,21:1]
42    [@41,2:5='0x00',<const_interger>,21:2]
43    [@42,7:7=';',<';'>,21:7]
44    [@43,1:8='0XABCDEF',<const_interger>,22:1]
45    [@44,10:10=';',<';'>,22:10]
46    [@45,1:5='0xful',<const_interger>,23:1]
47    [@46,7:7=';',<';'>,23:7]
48    [@47,1:6='0XFLLU',<const_interger>,24:1]
49    [@48,8:8=';',<';'>,24:8]
50    [@49,1:3='0.0',<const_float>,27:1]
51    [@50,5:5=';',<';'>,27:5]
52    [@51,1:1='+',<'+'>,28:1]
53    [@52,2:7='1.1e+1',<const_float>,28:2]
54    [@53,9:9=';',<';'>,28:9]
55    [@54,1:1='-',<'-'>,29:1]
56    [@55,2:7='1.1E-1',<const_float>,29:2]
57    [@56,9:9=';',<';'>,29:9]
58    [@57,1:6='1.1e1f',<const_float>,30:1]
59    [@58,8:8=';',<';'>,30:8]
60    [@59,1:6='1.1E1L',<const_float>,31:1]
61    [@60,8:8=';',<';'>,31:8]
62    [@61,1:5='0x0p0',<const_float>,33:1]
63    [@62,6:6=';',<';'>,33:6]
64    [@63,1:7='0x0.0p0',<const_float>,34:1]
65    [@64,9:9=';',<';'>,34:9]
66    [@65,1:1='+',<'+'>,35:1]
67    [@66,2:9='0xa.ap+1',<const_float>,35:2]
68    [@67,11:11=';',<';'>,35:11]
69    [@68,1:1='-',<'-'>,36:1]
70    [@69,2:9='0xa.aP-1',<const_float>,36:2]
71    [@70,11:11=';',<';'>,36:11]
72    [@71,1:8='0xa.ap1f',<const_float>,37:1]
73    [@72,10:10=';',<';'>,37:10]
74    [@73,1:8='0xa.aP1L',<const_float>,38:1]
75    [@74,10:10=';',<';'>,38:10]
76    [@75,1:3=''a'',<const_char>,41:1]
77    [@76,5:5=';',<';'>,41:5]
78    [@77,1:4='L'a'',<const_char>,42:1]
79    [@78,6:6=';',<';'>,42:6]
80    [@79,1:4='U'a'',<const_char>,43:1]
81    [@80,6:6=';',<';'>,43:6]
```

```
 82   [@81,1:4='u'a'',<const_char>,44:1]
 83   [@82,6:6=';',<';'>,44:6]
 84   [@83,1:4=''\n'',<const_char>,45:1]
 85   [@84,6:6=';',<';'>,45:6]
 86   [@85,1:4=''\?'',<const_char>,46:1]
 87   [@86,6:6=';',<';'>,46:6]
 88   [@87,1:5=''\24'',<const_char>,47:1]
 89   [@88,7:7=';',<';'>,47:7]
 90   [@89,1:17='"abcdefg123456\\"',<const_string>,50:1]
 91   [@90,19:19=';',<';'>,50:19]
 92   [@91,1:5='u8"a"',<const_string>,51:1]
 93   [@92,7:7=';',<';'>,51:7]
 94   [@93,1:4='u"a"',<const_string>,52:1]
 95   [@94,6:6=';',<';'>,52:6]
 96   [@95,1:4='U"a"',<const_string>,53:1]
 97   [@96,6:6=';',<';'>,53:6]
 98   [@97,1:4='L"a"',<const_string>,54:1]
 99   [@98,6:6=';',<';'>,54:6]
100   [@99,1:3='int',<keyword>,57:1]
101   [@100,5:7='__a',<Identifier>,57:5]
102   [@101,9:9=';',<';'>,57:9]
103   [@102,1:3='int',<keyword>,58:1]
104   [@103,5:7='__1',<Identifier>,58:5]
105   [@104,9:9=';',<';'>,58:9]
106   [@105,1:3='int',<keyword>,59:1]
107   [@106,5:6='a1',<Identifier>,59:5]
108   [@107,8:8=';',<';'>,59:8]
109   [@108,1:3='int',<keyword>,60:1]
110   [@109,5:5='a',<Identifier>,60:5]
111   [@110,7:7=';',<';'>,60:7]
112   [@111,1:3='int',<keyword>,61:1]
113   [@112,5:7='a_1',<Identifier>,61:5]
114   [@113,9:9=';',<';'>,61:9]
115   [@114,1:3='int',<keyword>,62:1]
116   [@115,5:6='a_',<Identifier>,62:5]
117   [@116,8:8=';',<';'>,62:8]
118   [@117,1:3='int',<keyword>,65:1]
119   [@118,5:5='i',<Identifier>,65:5]
120   [@119,6:6='=',<'='>,65:6]
121   [@120,7:7='1',<const_interger>,65:7]
122   [@121,9:9=';',<';'>,65:9]
123   [@122,1:5='float',<keyword>,66:1]
124   [@123,7:7='f',<Identifier>,66:7]
125   [@124,9:9=';',<';'>,66:9]
126   [@125,1:6='double',<keyword>,67:1]
127   [@126,8:8='d',<Identifier>,67:8]
128   [@127,10:10=';',<';'>,67:10]
129   [@128,1:4='char',<keyword>,68:1]
130   [@129,6:6='c',<Identifier>,68:6]
131   [@130,8:8=';',<';'>,68:8]
132   [@131,1:4='long',<keyword>,69:1]
133   [@132,6:6='l',<Identifier>,69:6]
```

```
134    [@133,8:8=';',<';'>,69:8]
135    [@134,1:5='short',<keyword>,70:1]
136    [@135,7:7='s',<Identifier>,70:7]
137    [@136,9:9=';',<';'>,70:9]
138    [@137,1:6='signed',<keyword>,71:1]
139    [@138,8:9='si',<Identifier>,71:8]
140    [@139,11:11=';',<';'>,71:11]
141    [@140,1:8='unsigned',<keyword>,72:1]
142    [@141,10:14='short',<keyword>,72:10]
143    [@142,16:17='us',<Identifier>,72:16]
144    [@143,19:19=';',<';'>,72:19]
145    [@144,1:7='typedef',<keyword>,73:1]
146    [@145,9:14='struct',<keyword>,73:9]
147    [@146,1:1='{',<'{'>,74:1]
148    [@147,2:4='int',<keyword>,75:2]
149    [@148,6:9='num1',<Identifier>,75:6]
150    [@149,10:10=';',<';'>,75:10]
151    [@150,1:1='}',<'}'>,76:1]
152    [@151,2:5='test',<Identifier>,76:2]
153    [@152,6:6=';',<';'>,76:6]
154    [@153,1:4='test',<Identifier>,77:1]
155    [@154,6:10='test1',<Identifier>,77:6]
156    [@155,12:12=';',<';'>,77:12]
157    [@156,1:4='test',<Identifier>,78:1]
158    [@157,6:6='*',<'*'>,78:6]
159    [@158,7:11='test2',<Identifier>,78:7]
160    [@159,13:13=';',<';'>,78:13]
161    [@160,1:6='static',<keyword>,79:1]
162    [@161,8:10='int',<keyword>,79:8]
163    [@162,12:14='sti',<Identifier>,79:12]
164    [@163,16:16=';',<';'>,79:16]
165    [@164,1:5='const',<keyword>,80:1]
166    [@165,7:9='int',<keyword>,80:7]
167    [@166,11:12='ci',<Identifier>,80:11]
168    [@167,14:14=';',<';'>,80:14]
169    [@168,1:6='sizeof',<keyword>,81:1]
170    [@169,7:7='(',<'('>,81:7]
171    [@170,9:11='int',<keyword>,81:9]
172    [@171,13:13=')',<')'>,81:13]
173    [@172,15:15=';',<';'>,81:15]
174    [@173,1:2='if',<keyword>,82:1]
175    [@174,3:3='(',<'('>,82:3]
176    [@175,5:5='1',<const_interger>,82:5]
177    [@176,7:7=')',<')'>,82:7]
178    [@177,1:1='{',<'{'>,83:1]
179    [@178,1:1='}',<'}'>,85:1]
180    [@179,1:4='else',<keyword>,86:1]
181    [@180,1:1='{',<'{'>,87:1]
182    [@181,1:1='}',<'}'>,89:1]
183    [@182,1:3='for',<keyword>,90:1]
184    [@183,4:4='(',<'('>,90:4]
185    [@184,6:6=';',<';'>,90:6]
```

```
186    [@185,8:8=';',<';'>,90:8]
187    [@186,10:10=')',<')'>,90:10]
188    [@187,1:1='{',<'{'>,91:1]
189    [@188,2:9='continue',<keyword>,92:2]
190    [@189,11:11=';',<';'>,92:11]
191    [@190,1:1='}',<'}'>,93:1]
192    [@191,1:5='while',<keyword>,94:1]
193    [@192,6:6='(',<'('>,94:6]
194    [@193,8:8='1',<const_interger>,94:8]
195    [@194,10:10=')',<')'>,94:10]
196    [@195,1:1='{',<'{'>,95:1]
197    [@196,1:1='}',<'}'>,97:1]
198    [@197,1:6='switch',<keyword>,98:1]
199    [@198,7:7='(',<'('>,98:7]
200    [@199,8:8='i',<Identifier>,98:8]
201    [@200,9:9=')',<')'>,98:9]
202    [@201,10:10='{',<'{'>,98:10]
203    [@202,5:8='case',<keyword>,99:5]
204    [@203,10:10='1',<const_interger>,99:10]
205    [@204,11:11=':',<':'>,99:11]
206    [@205,14:18='break',<keyword>,99:14]
207    [@206,19:19=';',<';'>,99:19]
208    [@207,5:11='default',<keyword>,100:5]
209    [@208,12:12=':',<':'>,100:12]
210    [@209,15:19='break',<keyword>,100:15]
211    [@210,20:20=';',<';'>,100:20]
212    [@211,1:1='}',<'}'>,101:1]
213    [@212,1:2='do',<keyword>,102:1]
214    [@213,1:1='{',<'{'>,103:1]
215    [@214,1:1='}',<'}'>,105:1]
216    [@215,2:6='while',<keyword>,105:2]
217    [@216,7:7='(',<'('>,105:7]
218    [@217,9:9='0',<const_interger>,105:9]
219    [@218,11:11=')',<')'>,105:11]
220    [@219,13:13=';',<';'>,105:13]
221    [@220,1:4='goto',<keyword>,106:1]
222    [@221,6:13='gotoFlag',<Identifier>,106:6]
223    [@222,15:15=';',<';'>,106:15]
224    [@223,1:3='int',<keyword>,108:1]
225    [@224,5:9='array',<Identifier>,108:5]
226    [@225,10:10='[',<'['>,108:10]
227    [@226,11:12='10',<const_interger>,108:11]
228    [@227,13:13=']',<']'>,108:13]
229    [@228,15:15='=',<'='>,108:15]
230    [@229,17:17='{',<'{'>,108:17]
231    [@230,18:18='0',<const_interger>,108:18]
232    [@231,19:19='}',<'}'>,108:19]
233    [@232,21:21=';',<';'>,108:21]
234    [@233,1:5='test1',<Identifier>,109:1]
235    [@234,6:6='.',<'.'>,109:6]
236    [@235,7:10='num1',<Identifier>,109:7]
237    [@236,12:12='=',<'='>,109:12]
```

```
238    [@237,14:14='0',<const_interger>,109:14]
239    [@238,16:16=';',<';'>,109:16]
240    [@239,1:5='test2',<Identifier>,110:1]
241    [@240,6:7='->',<'->'>,110:6]
242    [@241,8:11='num1',<Identifier>,110:8]
243    [@242,13:13='=',<'='>,110:13]
244    [@243,15:15='0',<const_interger>,110:15]
245    [@244,17:17=';',<';'>,110:17]
246    [@245,1:1='i',<Identifier>,111:1]
247    [@246,2:3='++',<'++'>,111:2]
248    [@247,5:5=';',<';'>,111:5]
249    [@248,1:1='i',<Identifier>,112:1]
250    [@249,2:3='--',<'--'>,112:2]
251    [@250,5:5=';',<';'>,112:5]
252    [@251,1:1='i',<Identifier>,113:1]
253    [@252,3:3='=',<'='>,113:3]
254    [@253,5:5='1',<const_interger>,113:5]
255    [@254,7:7='+',<'+'>,113:7]
256    [@255,9:9='1',<const_interger>,113:9]
257    [@256,11:11=';',<';'>,113:11]
258    [@257,1:1='i',<Identifier>,114:1]
259    [@258,3:3='=',<'='>,114:3]
260    [@259,5:5='1',<const_interger>,114:5]
261    [@260,7:7='-',<'-'>,114:7]
262    [@261,9:9='1',<const_interger>,114:9]
263    [@262,11:11=';',<';'>,114:11]
264    [@263,1:1='i',<Identifier>,115:1]
265    [@264,3:3='=',<'='>,115:3]
266    [@265,5:5='1',<const_interger>,115:5]
267    [@266,7:7='*',<'*'>,115:7]
268    [@267,9:9='1',<const_interger>,115:9]
269    [@268,11:11=';',<';'>,115:11]
270    [@269,1:1='i',<Identifier>,116:1]
271    [@270,3:3='=',<'='>,116:3]
272    [@271,5:5='1',<const_interger>,116:5]
273    [@272,7:7='/',<'/'>,116:7]
274    [@273,9:9='1',<const_interger>,116:9]
275    [@274,11:11=';',<';'>,116:11]
276    [@275,1:1='i',<Identifier>,117:1]
277    [@276,3:3='=',<'='>,117:3]
278    [@277,5:5='1',<const_interger>,117:5]
279    [@278,7:7='%',<'%'>,117:7]
280    [@279,9:9='1',<const_interger>,117:9]
281    [@280,11:11=';',<';'>,117:11]
282    [@281,1:1='i',<Identifier>,118:1]
283    [@282,3:3='=',<'='>,118:3]
284    [@283,5:5='!',<'!'>,118:5]
285    [@284,6:6='i',<Identifier>,118:6]
286    [@285,8:8=';',<';'>,118:8]
287    [@286,1:1='i',<Identifier>,119:1]
288    [@287,3:3='=',<'='>,119:3]
289    [@288,5:5='i',<Identifier>,119:5]
```

```
290    [@289,7:7='&',<'&'>,119:7]
291    [@290,9:9='1',<const_interger>,119:9]
292    [@291,11:11=';',<';'>,119:11]
293    [@292,1:1='i',<Identifier>,120:1]
294    [@293,3:3='=',<'='>,120:3]
295    [@294,5:5='i',<Identifier>,120:5]
296    [@295,7:7='|',<'|'>,120:7]
297    [@296,9:9='1',<const_interger>,120:9]
298    [@297,11:11=';',<';'>,120:11]
299    [@298,1:1='i',<Identifier>,121:1]
300    [@299,3:3='=',<'='>,121:3]
301    [@300,5:5='i',<Identifier>,121:5]
302    [@301,7:8='&&',<'&&'>,121:7]
303    [@302,10:10='1',<const_interger>,121:10]
304    [@303,12:12=';',<';'>,121:12]
305    [@304,1:1='i',<Identifier>,122:1]
306    [@305,3:3='=',<'='>,122:3]
307    [@306,5:5='i',<Identifier>,122:5]
308    [@307,7:8='||',<'||'>,122:7]
309    [@308,10:10='1',<const_interger>,122:10]
310    [@309,12:12=';',<';'>,122:12]
311    [@310,1:1='i',<Identifier>,123:1]
312    [@311,3:3='=',<'='>,123:3]
313    [@312,5:5='i',<Identifier>,123:5]
314    [@313,7:7='^',<'^'>,123:7]
315    [@314,9:9='1',<const_interger>,123:9]
316    [@315,11:11=';',<';'>,123:11]
317    [@316,1:1='i',<Identifier>,124:1]
318    [@317,3:3='=',<'='>,124:3]
319    [@318,5:5='i',<Identifier>,124:5]
320    [@319,7:8='>>',<'>>'>,124:7]
321    [@320,10:10='1',<const_interger>,124:10]
322    [@321,12:12=';',<';'>,124:12]
323    [@322,1:1='i',<Identifier>,125:1]
324    [@323,3:3='=',<'='>,125:3]
325    [@324,5:5='i',<Identifier>,125:5]
326    [@325,7:8='<<',<'<<'>,125:7]
327    [@326,10:10='1',<const_interger>,125:10]
328    [@327,12:12=';',<';'>,125:12]
329    [@328,1:1='i',<Identifier>,126:1]
330    [@329,3:3='=',<'='>,126:3]
331    [@330,5:5='i',<Identifier>,126:5]
332    [@331,7:7='>',<'>'>,126:7]
333    [@332,9:9='1',<const_interger>,126:9]
334    [@333,11:11='?',<'?'>,126:11]
335    [@334,13:13='1',<const_interger>,126:13]
336    [@335,15:15=':',<':'>,126:15]
337    [@336,17:17='2',<const_interger>,126:17]
338    [@337,19:19=';',<';'>,126:19]
339    [@338,1:1='i',<Identifier>,127:1]
340    [@339,3:4='+=',<'+='>,127:3]
341    [@340,6:6='1',<const_interger>,127:6]
```

```
342    [@341,8:8=';',<';'>,127:8]
343    [@342,1:1='i',<Identifier>,128:1]
344    [@343,3:4='-=',<'-='>,128:3]
345    [@344,6:6='1',<const_interger>,128:6]
346    [@345,8:8=';',<';'>,128:8]
347    [@346,1:1='i',<Identifier>,129:1]
348    [@347,3:4='*=',<'*='>,129:3]
349    [@348,6:6='1',<const_interger>,129:6]
350    [@349,8:8=';',<';'>,129:8]
351    [@350,1:1='i',<Identifier>,130:1]
352    [@351,3:3='/',<'/'>,130:3]
353    [@352,4:4='=',<'='>,130:4]
354    [@353,6:6='1',<const_interger>,130:6]
355    [@354,8:8=';',<';'>,130:8]
356    [@355,1:1='i',<Identifier>,131:1]
357    [@356,3:4='%=',<'%='>,131:3]
358    [@357,6:6='1',<const_interger>,131:6]
359    [@358,8:8=';',<';'>,131:8]
360    [@359,1:1='i',<Identifier>,132:1]
361    [@360,3:4='&=',<'&='>,132:3]
362    [@361,6:6='1',<const_interger>,132:6]
363    [@362,8:8=';',<';'>,132:8]
364    [@363,1:1='i',<Identifier>,133:1]
365    [@364,3:4='|=',<'|='>,133:3]
366    [@365,6:6='1',<const_interger>,133:6]
367    [@366,8:8=';',<';'>,133:8]
368    [@367,1:1='i',<Identifier>,134:1]
369    [@368,3:4='^=',<'^='>,134:3]
370    [@369,6:6='1',<const_interger>,134:6]
371    [@370,8:8=';',<';'>,134:8]
372    [@371,1:1='i',<Identifier>,135:1]
373    [@372,3:5='>>=',<'>>='>,135:3]
374    [@373,7:7='1',<const_interger>,135:7]
375    [@374,9:9=';',<';'>,135:9]
376    [@375,1:1='i',<Identifier>,136:1]
377    [@376,3:5='<<=',<'<<='>,136:3]
378    [@377,7:7='1',<const_interger>,136:7]
379    [@378,9:9=';',<';'>,136:9]
380    [@379,1:2='if',<keyword>,137:1]
381    [@380,3:3='(',<'('>,137:3]
382    [@381,5:5='i',<Identifier>,137:5]
383    [@382,6:7='==',<'=='>,137:6]
384    [@383,8:8='1',<const_interger>,137:8]
385    [@384,10:10=')',<')'>,137:10]
386    [@385,11:11='{',<'{'>,137:11]
387    [@386,12:12='}',<'}'>,137:12]
388    [@387,1:2='if',<keyword>,138:1]
389    [@388,3:3='(',<'('>,138:3]
390    [@389,5:5='i',<Identifier>,138:5]
391    [@390,6:7='!=',<'!='>,138:6]
392    [@391,8:8='1',<const_interger>,138:8]
393    [@392,10:10=')',<')'>,138:10]
```

```
394    [@393,11:11='{',<'{'>,138:11]
395    [@394,12:12='}',<'}'>,138:12]
396    [@395,1:2='if',<keyword>,139:1]
397    [@396,3:3='(',<'('>,139:3]
398    [@397,5:5='i',<Identifier>,139:5]
399    [@398,6:6='>',<'>'>,139:6]
400    [@399,7:7='1',<const_interger>,139:7]
401    [@400,9:9=')',<')'>,139:9]
402    [@401,10:10='{',<'{'>,139:10]
403    [@402,11:11='}',<'}'>,139:11]
404    [@403,1:2='if',<keyword>,140:1]
405    [@404,3:3='(',<'('>,140:3]
406    [@405,5:5='i',<Identifier>,140:5]
407    [@406,6:6='<',<'<'>,140:6]
408    [@407,7:7='1',<const_interger>,140:7]
409    [@408,9:9=')',<')'>,140:9]
410    [@409,10:10='{',<'{'>,140:10]
411    [@410,11:11='}',<'}'>,140:11]
412    [@411,1:2='if',<keyword>,141:1]
413    [@412,3:3='(',<'('>,141:3]
414    [@413,5:5='i',<Identifier>,141:5]
415    [@414,6:7='>=',<'>='>,141:6]
416    [@415,8:8='1',<const_interger>,141:8]
417    [@416,10:10=')',<')'>,141:10]
418    [@417,11:11='{',<'{'>,141:11]
419    [@418,12:12='}',<'}'>,141:12]
420    [@419,1:2='if',<keyword>,142:1]
421    [@420,3:3='(',<'('>,142:3]
422    [@421,5:5='i',<Identifier>,142:5]
423    [@422,6:7='<=',<'<='>,142:6]
424    [@423,8:8='1',<const_interger>,142:8]
425    [@424,10:10=')',<')'>,142:10]
426    [@425,11:11='{',<'{'>,142:11]
427    [@426,12:12='}',<'}'>,142:12]
428    [@427,0:7='gotoFlag',<Identifier>,143:0]
429    [@428,8:8=':',<':'>,143:8]
430    [@429,1:6='return',<keyword>,144:1]
431    [@430,8:8='0',<const_interger>,144:8]
432    [@431,10:10=';',<';'>,144:10]
433    [@432,0:0='}',<'}'>,145:0]
434    [@433,0:3='void',<keyword>,146:0]
435    [@434,5:12='function',<Identifier>,146:5]
436    [@435,13:13='(',<'('>,146:13]
437    [@436,15:17='int',<keyword>,146:15]
438    [@437,19:22='arg1',<Identifier>,146:19]
439    [@438,24:24=',',<','>,146:24]
440    [@439,26:28='int',<keyword>,146:26]
441    [@440,30:33='arg2',<Identifier>,146:30]
442    [@441,35:35=')',<')'>,146:35]
443    [@442,0:0='{',<'{'>,147:0]
444    [@443,0:0='}',<'}'>,148:0]
```

上图显示，识别出了444个单词，识别基本正确。

## 6. 实验感想

在这次的实验过程中，我直接利用正则表达式进行书写Flex的.l文件生成，一方面我是对正则表达式用了更深的理解和运用，另一方面，我对C语言的规范也有了进一步的认识，虽然某些规则看起来很简单，但是在书写的时候要考虑到各种细节。不足的是，虽然利用Flex可以很好的匹配相应类别的单词，但缺少了错误处理的过程，在词法分析阶段，有些错误可以处理的可能要推迟到语法分析阶段，这样以来，编译的性能可能会有所减弱。

除此之外，我也理解了计卫星老师所给的BIT-minCC的框架，学会了如何在里面嵌入自己的代码和可执行程序，对整个编译过程有了更为熟悉的认识，要较好的实现编译器的前端，并不是一件容易的事。

对词法分析的设计与实现有了深刻的体会，尤其DFA的设计层面，如何完成单词的识别，如何对整个DFA进行汇总，刚开始我设计的DFA有很多缺陷不足，需要不断地修改，尽管现在也不是很完美，但是基本的识别功能已经完成。