

实验题目

词法分析扫描器的设计

实验要求

- 在报告中明确词法分析器所支持的单词范围
- 有自动机的设计
- 有实验结果说明

实验过程

此处实现的是 C 语言的词法分析器。

(1) 首先定义标识符、保留字、常数、运算符、界符和需略除的符号的范围。

标识符	标识符可以是字母、数字、下划线(A~Z, a~z, 0~9, _)组成的字符串，并且第一个字符必须是字母或下划线。
保留字	指在高级语言中已经定义过的字，使用者不能再将这些字作为变量名或过程名使用。
常数	数字0~9
运算符	包括算术运算符、关系运算符、逻辑运算符、赋值运算符、位运算符、其他运算符
界符	限制界限所用，有单字节界符和双字节界符
需略除的符号	空格，制表符，换行符

标识符：

1. 标识符是严格区分大小写的。
2. 标识符**不能是C语言的关键字和保留标识符**。
3. 标识符长度限制跟C语言标准和编译器环境有关。

保留字：

字	int,long,short,float,double,char,unsigned,signed,const,void,volatile,enum,struct,union
语句定义保留字	if,else,goto,switch,case,do,while,for,continue,break,return,default,typedef
存储类说明保留字	auto,register,extern,static
长度运算符保留字	sizeof

运算符：

算术运算符	+, -, *, /, %, ++, --
关系运算符	==, !=, >, <, >=, <=
逻辑运算符	&&, , !
赋值运算符	=, +=, -=, *=, /=, %=, <<=, >>=, &=, ^=, =
位运算符	&, , ^, ~, <<, >>
其他运算符	Condition ? X:Y, ., ->, &, *, ,

界符：

单字节界符	;
双字节界符	{, }, [,], " 和 ", ' 和 '

需略除的符号：

注释符	\\, /* 和 */,
其他	空格, 制表符, 换行符

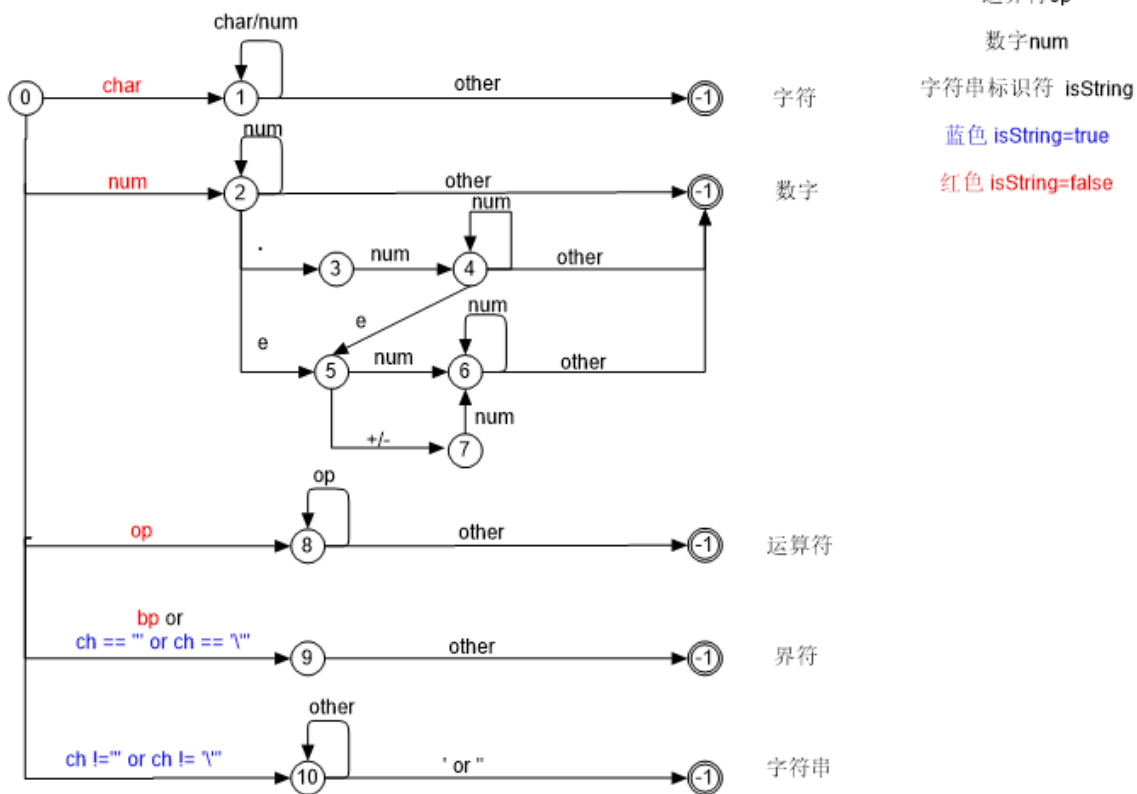
(2) 进行种别编码：

单词符号	种别编码	助记符	内码值
常数	0	num	num在常数表中的位置
标识符	1	id	id在符号表中的位置
字符串	2	s	s在字符串表中的位置
int	3	int	
long	4	long	
short	5	short	
float	6	float	
double	7	double	
char	8	char	
unsigned	9	unsigned	
signed	10	signed	
const	11	const	
void	12	void	
volatile	13	volatile	
enum	14	enum	
struct	15	struct	
union	16	union	
if	17	if	-
else	18	else	-
goto	19	goto	-
switch	20	switch	-
case	21	case	-
do	22	do	-
while	23	while	-
for	24	for	-
continue	25	continue	-
break	26	break	-
return	27	return	-
default	28	default	-
typedef	29	typedef	-

单词符号	种别编码	助记符	内码值
auto	30	auto	-
register	31	register	-
extern	32	extern	-
static	33	static	-
sizeof	34	sizeof	-
+	35	+	-
-	36	-	-
*	37	*	-
/	38	/	-
%	39	%	-
++	40	++	-
--	41	--	-
==	42	==	-
!=	43	!=	-
>	44	>	-
<	45	<	-
>=	46	>=	-
<=	47	<=	-
&&	48	&&	-
	49		-
!	50	!	-
=	51	=	-
+=	52	+=	-
-=	53	-=	-
*=	54	*=	-
/=	55	/=	-
%=	56	%=	-
<<=	57	<<=	-
>>=	58	>>=	-
&=	59	&=	-

单词符号	种别编码	助记符	内码值
^=	60	^=	-
=	61	=	-
&	62	&	-
	63		-
^	64	^	-
~	65	-	-
<<	66	<<	-
>>	67	>>	-
;	68	;	-
{	69	{	-
}	70	}	-
(71	(-
)	72)	-
[73	[-
]	74]	-
"	75	"	-
'	76	'	-
,	77	,	-
#	78	#	-
.	79	.	-

(3) 建立自动机



(4) 代码实现

4.1 用数组存储关键字、运算符、界符等常量信息。

```

1  const char *keyword[] = //定义关键字
2  {"int","long","short","float","double","char","unsigned",
3  "signed","const","void","volatile","enum","struct","union",
4  "if", "else", "goto", "switch", "case", "do", "while",
5  "for", "continue", "break", "return", "default",
6  "typedef", "auto", "register", "extern", "static", "sizeof"};
7  int keyword_length = sizeof(keyword) / sizeof(char*);
8  const char *Opera[] = //定义运算符
9  { "+", "-", "*", "/", "%", "++", "--", "==", "!=", ">",
10  "<", ">=", "<=", "&&", "||", "!", "=", "+=", "-=", "*=",
11  "/=", "%=", "<<=", ">>=", "&=", "^=", "|=", "&", "|", "^",
12  "~", "<<", ">>" };
13  const char *BS[] = { ";", "{", "}", "(", ")", "[", "]", "\\", " ", " ", " ", "#", "." }; //
    定义界符
  
```

4.2 函数和变量说明

```

1  ifstream inFile; //需要编译的文件
2  int BS_length = sizeof( BS ) / sizeof(char*);
3  int operator_length = sizeof(Opera) / sizeof(char*);
4  char id_table[255][255]; //标识符表，存放构成的标识符
5  int ptr_id = 0; //标识符表的下标
6  char num_table[255][255]; //常数表，存放构成的常熟
7  int ptr_num = 0; //常数表的下标
8  char string_table[255][255];
9  int ptr_string = 0;
  
```

```

10 char ch; //字符变量，存放最新读入的源程序字符
11 int ptr_token = 0; //token数组的下标
12 char token[255] = ""; //存放当前构成单词符号的字符串
13 bool iserror = false; // 发生错误的标志
14 bool isdigit(); //判断当前读取的字符是否是数字
15 bool isletter(); //判断当前读取的字符是否是字母
16 bool isoperator(); //判断当前字符是否是运算符
17 bool isbp(); //判断当前字符是否是界符
18 void getbe(); //过滤无用符号
19 void concatenation(); //将ch中的字符链接到token后作为新的字符串
20 void error(int num); // 发现错误，显示错误
21 int reserve(); //判断token是否是keyword，若是则返回它的种别编码，若不是返回-1
22 int opara(); //判断token是否是操作符，若是，则返回它的种别编码，若不是返回-1
23 int bp(); //判断token是否是界符，若是，则返回它的种别编码，若不是返回-1
24 void buildtable(); //建表（标识符表或常数表）
25 void automata(); //自动机
26 void LexicalAnalyzer(char*); //整个词法分析器
27 void Getchar();
28 int isString = 0; //用于标记token是字符串还是标识符，根据前一个字符是否是'或"决定，若是字符串isString=1,否则为0
29 char pre; //保存上一个字符
30 bool useNext = false;
31 char nextch;

```

4.3 函数实现

- bool isdigit()

```

1 bool isdigit(){
2     if( ch >= '0' && ch <= '9' )
3         return true;
4     else
5         return false;
6 }

```

此函数用于判断当前字符是否是数字 0 ~ 9。

- bool isletter()

```

1 bool isletter(){
2     if( (ch >= 'a' && ch <= 'z') || (ch >= 'A' && ch <= 'Z') )
3         return true;
4     else
5         return false;
6 }

```

此函数用于判断当前字符是否是字母。

- bool isoperator()

```

1 bool isoperator(){
2     if( ch == '+' || ch == '-' || ch == '*' || ch == '/' || ch == '%' ||
3     ch == '=' || ch == '!' || ch == '>' || ch == '<' || ch == '&' ||
4     ch == '|' || ch == '^' || ch == '~')
5         return true;
6     return false;
7 }

```

此函数用于判断当前字符是否是操作符，直接枚举判断。

- `bool isbp()`

```

1 bool isbp(){
2     if( ch == ';' || ch == '{' || ch == '}' || ch == '(' || ch == ')' || ch ==
3     '['
4     || ch == ']' || ch == '"' || ch == '\'' || ch == ',' || ch == '#' || ch ==
5     '.' )
6         return true;
7     return false;
8 }

```

此函数用于判断当前字符是否是界符，直接枚举判断。

- `void getbe()`

```

1 void getbe(){
2     while(ch == ' ' || ch == '\t' || ch == '\n')
3         Getchar();
4 }

```

此函数用于过滤无用字符。其中定义无用字符为空格，制表符和换行符。

- `void error(int num)`

```

1 void error(int num){
2     if( num == 0 )
3         cout << "The token is too long\n" ;
4     if( num == 1 )
5         cout << "Illegal character: " << ch << endl;
6 }

```

此函数用于函数执行到错误处时报错，并输出错误类型的日志信息。传进的参数 `num` 为 error 的编码。这里定义了两种错误，第一种是 token 的长度超出最大限制，这里限制 token 最长为 254 个字符；第二种是接收到了非法字符。

- `void concatenation()`


```

1 void concatenation(){
2     if (ptr_token > 254){
3         error(0);
4         return;
5     }
6     token[ptr_token] = ch;
7     ptr_token ++;
8     token[ptr_token] = '\0';
9 }

```

此函数将当前字符加入 `token`，如果 `token` 的长度大于 254，则返回错误。否则将当前字符加入 `token` 尾端。

- `int reserve()`

```

1 int reserve(){
2     int i = 0;
3     while( i < keyword_length ){
4         if( !strcmp( keyword[i], token) )
5             return i + 3; //返回种别编码
6         i++;
7     }
8     return -1; //不是关键字，返回-1
9 }

```

此函数用于判断 `token` 是否是关键字，若是则返回它的种别编码，若不是返回-1。判断方法是將 `token` 与关键字数组中的每个元素进行对比，判断它们是否相等，如果相等说明 `token` 是关键字。

- `int opra()`

```

1 int opra(){
2     int i = 0;
3     while( i < operator_length ){
4         if( !strcmp(Opera[i], token) )
5             return i + 35;
6         i++;
7     }
8     if( ptr_token <= 4 && ptr_token > 2) // 有可能是两个操作符连在一起，判断不出，
    将其拆分后重新进行判断。
9     {
10         if( pre == '+' || pre == '-' || pre == '*' || pre == '/' || pre ==
    '%' ||
11         pre == '=' || pre == '!' || pre == '>' || pre == '<' || pre == '&' ||
12         pre == '|' || pre == '^' || pre == '~' ){
13             ch = pre; //开始回滚
14             pre = token[ptr_token - 2];
15             token[ptr_token - 1] = '\0';
16             useNext = true; // 回滚，char下次更新时使用next的值，而不直接读取文件
17             i = 0;
18             while( i < operator_length ){
19                 if( !strcmp(Opera[i], token) ){
20                     return i + 35;
21                 }
22                 i++;
23             }
24         }
    }

```

```

25     }
26     return -1;
27 }

```

此函数用于判断 `token` 是否是操作符，若是则返回它的种别编码，若不是返回-1。判断方法是 将 `token` 与操作符数组中的每个元素进行对比，判断它们是否相等，如果相等说明 `token` 是操作符。返回的编码是 `i + 35`，原因是第一个操作符在上述表格中的种别编码为 35。另外，根据上述自动机的逻辑，当出现两个操作符连在一起，又不是合法操作符时，可能无法正确判别，此时需要进行回滚操作。

- `int bp()`

```

1  int bp(){
2      int i = 0;
3      while( i < BS_length ){
4          if( !strcmp(BS[i], token) )
5              return i + 68;
6          i++;
7      }
8      return -1;
9  }

```

该函数用于判断 `token` 是否是界符，若是，则返回它的种别编码，若不是返回-1。返回的编码是 `i + 68`，原因是第一个界符在上述表格中的种别编码为 68。

- `void Getchar()`

```

1  void Getchar(){
2      if (useNext){
3          ch = nextCh;
4          useNext = false;
5      }
6      else{
7          inFile >> ch;
8          nextCh = ch;
9      }
10 }

```

该函数用于读取字符。用一个布尔变量 `useNext` 判断是否要回滚，如果 `useNext` 为真，则不读取文件，而是用 `nextCh` 的值。`nextCh` 的值保存的是最近一次读取文件的字符。

- `void buildtable()`

```

1  void buildtable(){
2      if(isString != 0){
3          int k = 0;
4          while( k < ptr_string){
5              if( !strcmp( string_table[k], token ) ) {
6                  break;
7              }
8              k++;
9          }
10         if( k == ptr_string ){
11             strcpy(string_table[k], token);
12             ptr_string ++;
13         }
14         //isString = 0;

```

```

15     }
16     else{
17         char head = token[0]-'0';
18         if(!(head >= 0 && head <= 9)){ //token是标识符
19             int i = 0;
20             while( i < ptr_id ){
21                 if( !strcmp(id_table[i], token) ){ //两个字符串相等
22                     break;
23                 }
24                 i++;
25             }
26             if( i == ptr_id ){
27                 strcpy(id_table[ptr_id], token);
28                 ptr_id++;
29             }
30         }
31
32         else{ //token是常数
33             int j = 0;
34             while( j < ptr_num ){
35                 if( !strcmp(num_table[j], token) ){ //两个字符串相等
36                     break;
37                 }
38                 j ++;
39             }
40             if( j == ptr_num ){
41                 strcpy(num_table[ptr_num], token);
42                 ptr_num ++;
43             }
44         }
45     }
46 }

```

该函数用于建表（标识符表或常数表），为被判定为标识符或常数的 `token` 建表，保存出现过的标识符和常数。如果当前 `token` 被判断为标识符（常数），则在标识符表（常数表）逐个判断 `token` 是否已经存在于标识符表（常数表），若存在，则结束。否则，在标识符表（常数表）中保存该标识符（常数）。

- `void automata()`

```

1 void automata(){
2     int state = 0;
3     while(state != -1){ //规定-1为退出状态
4         if(state == 0){
5             if(ch == '"' || ch == '\\' )//不能直接isString == 1就跳到10
6                 state = 9;
7             else if( isString == 1 ){
8                 state = 10;
9             }
10            else{
11                getbe();
12                if(isletter())
13                    state = 1;
14                else if( isdigit() )
15                    state = 2;
16                else if( isoperator() )
17                    state = 8;

```

```

18         else if( isbp() )
19             state = 9;
20         else{
21             error(1);
22             cout << "191" << endl;
23             iserror = true;
24         }
25     }
26 }
27
28 else if( state == 1 ){
29     if( isletter() || isdigit() )
30         state = 1;
31     else
32         state = -1;
33 }
34 else if( state == 2 ){
35     if( isdigit() )
36         state = 2;
37     else if( ch == '.' )
38         state = 3;
39     else if( ch == 'e' )
40         state = 5;
41     else
42         state = -1;
43 }
44 else if( state == 3 ){
45     if( isdigit() )
46         state = 4;
47     else{
48         error(1);
49         cout << "217" << endl;
50         iserror = true;
51     }
52 }
53 else if( state == 4 ){
54     if( isdigit() )
55         state = 4;
56     else if( ch == 'e' )
57         state = 5;
58     else
59         state = -1;
60 }
61 else if( state == 5 ){
62     if( isdigit() ){
63         state = 6;
64     }
65     else if( ch == '+' || ch == '-' ){
66         state = 7;
67     }
68     else{
69         error(1);
70         cout << "238" << endl;
71         iserror = true;
72     }
73 }
74 else if( state == 6 ){
75     if( isdigit() ){

```

```

76         state = 6;
77     }
78     else
79         state = -1;
80 }
81 else if( state == 7){
82     if( isdigit() )
83         state = 6;
84     else {
85         error(1);
86         cout << "254" << endl;
87         iserror = true;
88     }
89 }
90 else if( state == 8){
91     if( isoperator() )
92         state = 8;
93     else
94         state = -1;
95 }
96 else if( state == 9 ){
97     state = -1;
98 }
99 else if( state == 10 ){
100     if( ch == '"' || ch == '\'' ){
101         state = -1;
102     }
103 }
104 if( iserror ) {
105     cout << "error happened in automata" << endl;
106     getchar();
107     //exit(EXIT_FAILURE);
108 }
109 if(state != -1){
110     concatenation();
111     pre = ch;
112     Getchar();
113 }
114 }
115 return;
116 }

```

该函数为自动机的实现，逻辑见（3）中的图。代码中，用一个变量 `state` 表示当前状态，然后根据当前字符，判断应该进入哪个状态。终止状态为 -1。

- `void LexicalAnalyzer(char* fileName)`

```

1 void LexicalAnalyzer(char* fileName){
2     ofstream outFile;
3     inFile.open(fileName, ios::in);
4     inFile.unsetf(ios_base::skipws); //逐字符读取文件，不跳过空白符
5     outFile.open("./LexicalAnalyzerOut.txt", ios::out);
6     if( !inFile.is_open() ){
7         cout << "Could not find the file!\n";
8         getchar();
9         //exit(EXIT_FAILURE);
10    }

```

```

11 Getchar();
12 outFile.setf(ios::left);
13 outFile << setfill(' ')<<setw(66) << "Characters";
14 outFile << "Code" << endl;
15 while(!inFile.eof()){
16     ptr_token = 0;
17     automata();
18     int flag = -1;
19     if( !( pre == '"' || pre == '\\' ) && isString == 1 ){ //token是字符
串
20         flag = 2;
21         outFile.setf(ios::left);
22         outFile << setfill(' ')<<setw(75) << token;
23         outFile << flag << endl;
24         //isString = 0;
25     }
26     else{
27         flag = reserve();
28         if( flag != -1 ){ //token为关键字
29             outFile.setf(ios::left);
30             outFile << setfill(' ')<<setw(75) << token ;
31             outFile << flag << endl;
32         }
33         else{
34             flag = opra();
35             if( flag != -1 ){
36                 outFile.setf(ios::left);
37                 outFile << setfill(' ')<<setw(75) << token;
38                 outFile << flag << endl;
39             }
40             else{
41                 flag = bp();
42                 if( flag != -1 ){ //token是界符
43                     outFile.setf(ios::left);
44                     outFile << setfill(' ')<<setw(75) << token;
45                     outFile << flag << endl;
46                     buildtable();
47                     if( pre == '"' || pre == '\\' )
48                         isString = isString==1 ? 0 : 1;
49                 }
50             }
51             else{ //token是标识符或常数
52                 if( token[0] >= '0' && token[0] <= '9' ) { //token是
常数
53                     flag = 0;
54                     outFile.setf(ios::left);
55                     outFile << setfill(' ')<<setw(75) << token;
56                     outFile << flag << endl;
57                     buildtable();
58                 }
59                 else{//token是标识符
60                     flag = 1;
61                     outFile.setf(ios::left);
62                     outFile << setfill(' ')<<setw(75) << token;
63                     outFile << flag << endl;
64                     buildtable();
65                 }
66             }

```

```

67         }
68     }
69 }
70 }
71 outFile.close();
72 inFile.close();
73 }

```

该函数实现整个词法分析器的功能。首先是打开要读取的文件，读取文件字符，调用自动机实现函数，得到 `token`，再调用判断 `token` 的函数，得到该 `token` 的种别编码信息。如果到此还没有报错，则将 `token` 和其种别编码信息写入文件 `./LexicalAnalyzerOut.txt`。循环读取文件字符，直到字符串读取完毕。最后关闭打开的两个文件。

- `int main()`

```

1  int main(){
2      char* fileName = "test.c";
3      LexicalAnalyzer(fileName);
4      cout << "Finished\n";
5      getchar();
6      return 0;
7  }

```

在 `main` 函数中指定要读取的文件名字，分析完毕后输出提示信息 "Finished"。

至此，我的词法分析器已实现完毕。

实验测试

测试代码为在网上 copy 的一段实现汉诺塔的 C 语言代码，如下：

```

1  #include <stdio.h>
2  int main()
3  {
4      int hanoi(int,char,char,char);
5      int n,counter;
6      printf("Input the number of disk: ");
7      scanf("%d",&n);
8      printf("\n");
9      counter=hanoi(n,'A','B','C');
10     return 0;
11 }
12 int hanoi(int n,char x,char y,char z)
13 {
14     int move(char,int,char);
15     if(n==1)
16         move(x,1,z);
17     else
18     {
19         hanoi(n-1,x,z,y);
20         move(x,n,z);
21         hanoi(n-1,y,x,z);
22     }
23     return 0;
24 }

```

```

25  int move(char getone,int n,char putone)
26  {
27      static int k=1;
28      printf("%2d:%3d # %c---%c\n",k,n,getone,putone);
29      if(k++%3==0)
30          printf("\n");
31      return 0;
32  }

```

实验结果

LexicalAnalyzerOut.txt - 记事本

文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

Characters	Code
#	78
include	1
<	45
stdio	1
.	79
h	1
>	44
int	3
main	1
(71
)	72
{	69
int	3
hanoi	1
(71
int	3
,	77
char	8
,	77
char	8
,	77
char	8
)	72
;	68
int	3
n	1
,	77
counter	1

第 1 行, 第 1 列 100% Windows (CRLF) UTF-8

实验输出文件 LexicalAnalyzerOut.txt 将打包在压缩文件中。

实验感想

(1) 此次实验让我对词法分析器的作用有了深刻的了解。

词法分析器的主要任务是读入源程序的输入字符、将它们组成词素，生成并输出一个词法单元序列，每个词法单元对应于一个词素。

当词法分析器发现了一个标识符的词素时，要将这个词素添加到符号表中。

同时词法分析器还需要完成：

- 过滤掉源程序中的注释和空白。
- 将编译器生成的错误信息与源程序的位置联系起来。记录行号等。

(2) 在设计词法分析器时，先要确定该语言的保留字、支持的运算符以及界符有哪些，然后再确定你要实现的词法分析器能支持到什么程度，然后进一步设计自动机，和进行代码的实现。

(3) 自动机的设计体现了词法分析器的运行逻辑，自动机设计得好可以让代码更简洁明了。

(4) 此处实现的功能有

- 对C语言全部关键字、运算符、界符进行识别；
- 对常数、字符串和标识符进行识别和存储；
- 对字符串的合法性进行判断。