

# 编译原理第一次实验

---

学号：19335299

姓名：朱德朋

## 实验题目

---

词法分析扫描器的设计实现

## 实验目的

---

了解高级语言单词的分类；

了解状态图以及如何表示并识别单词规则，掌握状态图到识别程序的编程；

## 实验内容

---

### 使用LEX(FLEX)产生词法分析器的要求

用LEX生成XX(以C为例)语言的词法分析器

#### 词法规则

- 了解所选择编程语言单词符号及其种别值

#### 功能

- 输入一个C语言源程序文件demo.c
- 输出一个文件tokens.txt，该文件包括每一个单词及其种类枚举值，每行一个单词

#### 提交6个文件

- 实验报告
- C语言的LEX源程序：clang.lex
- C语言词法分析程序C源程序：lex.yy.c
- C语言词法分析程序的可执行文件：clang.out/clang.exe
- C语言源程序文件：demo.c（实验输入）
- 词法分析及结果文件：tokens.txt（实验输出）

同时上传源码至Github

## 知识回顾

---

## 词法分析器的作用

**词法分析是编译的第一阶段。**词法分析器的功能输入源程序，按照构词规则分解成一系列单词符号。单词是语言中具有独立意义的最小单位，包括关键字、标识符、运算符、界符和常量等

1. 关键字 是由程序语言定义的具有固定意义的标识符。例如，Pascal 中的begin, end, if, while都是保留字。这些字通常不用作一般标识符。
2. 标识符 用来表示各种名字，如变量名，数组名，过程名等等。
3. 常数 常数的类型一般有整型、实型、布尔型、文字型等。
4. 运算符 如+、-、\*、/等等。
5. 界符 如逗号、分号、括号、等等。

把词法分析安排为一个独立阶段的好处是，它可以使整个编译程序的结构更简洁、清晰和条理化。

## 正规表达式与有限自动机

对于字母表 $\Sigma$ ，我们感兴趣的是它的一些特殊字集，即所谓正规集。我们将使用正规式这个概念来表示正规集。下面是正规式和正规集的递归定义：

1.  $\epsilon$ 和 $\phi$ 都是 $\Sigma$ 上的正规式，它们所表示的正规集分别为 $\{\epsilon\}$ 和 $\phi$ ；
2. 任何 $a \in \Sigma$ ， $a$ 是 $\Sigma$ 上的一个正规式，它所表示的正规集为 $\{a\}$ ；
3. 假定 $U$ 和 $V$ 都是 $\Sigma$ 上的正规式，它们所表示的正规集分别为 $L(U)$ 和 $L(V)$ ，那么， $(U)$ 、 $(U|V)$ 、 $(U \cdot V)$ 和 $U^*$ 也都是正规式，它们所表示的正规集分别为 $L(U) \cup L(V)$ 、 $L(U)L(V)$ 和 $(L(U))^*$ ；

仅由有限次使用上述三步骤而得到的表达式才是 $\Sigma$ 上的正规式。仅由这些正规式所表示的字集才是 $\Sigma$ 上的正规集。

正规式与有限自动机的等价性定理如下：

**对任何有限自动机  $M$ ，都存在一个正规式  $r$ ，使得  $L(r) = L(M)$ ；**

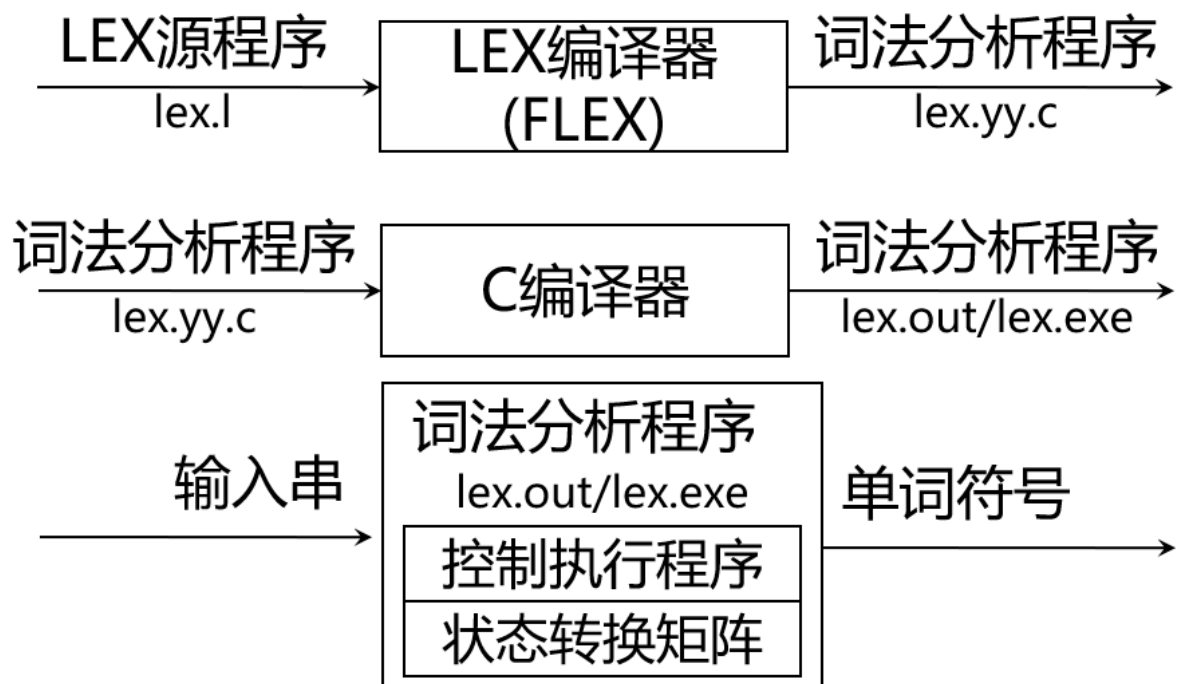
**对任何正规式  $r$ ，都存在一个有限自动机  $M$ ，使得  $L(M) = L(r)$ ；**

对于正规式 $r$ 转换成确定有限自动机(DFA)大致分为下面几步：

1. 正规表达式 $r$ 转化为非确定有限自动机(NFA)；
2. 非确定有限自动机(NFA)进一步确定化为确定有限自动机(DFA)；
3. 确定有限自动机(DFA)经过化简，得到最小化的确定有限自动机(DFA)；

## 词法分析器的自动产生

一个描述词法分析器的 LEX 程序由一组正规式以及与每个正规式相应的一个“动作”组成。“动作”本身是一小段程序代码，它指出了当按正规式识别出一个单词符号时应采取的行动。将 LEX 程序被编译后所得的结果程序记为 $L$ ，其作用就如同一个有限自动机一样，可用来识别和产生单词符号。结果程序含有一张状态转换表和一个控制程序。



一个LEX源程序主要包括两个部分，一部分是正规定义式，另一部分是识别规则。

如果  $\Sigma$  是一个字母表， $\Sigma$  上的正规定义式是下述形式的定义序列：

$$\begin{aligned}
 d_1 &\rightarrow r_1 \\
 d_2 &\rightarrow r_2 \\
 &\dots \\
 d_i &\rightarrow r_i
 \end{aligned}$$

其中  $d_i$  表示不同的名字，每个  $r_i$  是  $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$  上的符号所构成的正规式。 $r_i$  中不能含有  $d_{i+1}, d_{i+2}, \dots, d_n$ ，这样，对任何  $r_i$ ，可以构成一个  $\Sigma$  上的正规表达式，只要反复地将式中出现的名字代之以相应的正规式即可。注意，如果允许  $r_i$  中出现某些  $d_j$ ， $j \geq i$ ，那么这种替代过程将有可能不终止。

LEX 源程序中的识别规则是一串如下形式的LEX语句：

$$\begin{aligned}
 P_1 &\quad \{A_1\} \\
 P_2 &\quad \{A_2\} \\
 &\dots \\
 P_m &\quad \{A_m\}
 \end{aligned}$$

其中，每个  $P_i$  是一个正规式，称为词形。 $P_i$  中除了出现  $\Sigma$  中的字符外，还可以出现正规定义式左部所定义的任何简名  $d_i$ 。即， $P_i$  是  $\Sigma \cup \{d_1, d_2, \dots, d_n\}$  上的一个正规式。由于每个  $d_i$  最终都可化为纯粹  $\Sigma$  上的正规式，因此，每个  $P_i$  也同样如此，每个  $A_i$  是一小段程序代码，它指出了，在识别出词形为  $P_i$  的单词之后，词法分析器应采取的动作。这些识别规则完全决定了词法分析器 L 的功能。分析器 L 只能识别具有词形  $P_1, \dots, P_m$  的单词符号。

至于 LEX 所产生的目标程序 L（词法分析器）的工作原理：L 逐一扫描输入串的每个字符，寻找一个最长的子串匹配某个  $P_i$ ，将该子串截下来放在一个叫做 TOKEN 的缓冲区中（事实上，这个 TOKEN 也可以只包含一对指示器，它们分别指出这个子串在原输入缓冲区中的始末位置）。然后，L 就调用动作子程序  $A_i$ ，当  $A_i$  工作完后，L 就把所得的单词符号（由种别编码和属性值两部分构成）交给语法分析程序。当 L 重新被调用时就从输入串中继上次截出的位置之后识别下一个单词符号。

# 实验过程

根据上面知识回顾，我们知道：LEX的基本工作原理为：由正规式生成NFA，将NFA变换成DFA，DFA经化简后，模拟生成词法分析器。

其中正规式由开发者使用LEX语言编写，其余部分由LEX翻译器完成。翻译器将LEX源程序翻译成一个名为 `lex.yy.c` 的C语言源文件，此文件含有两部分内容：一部分是根据正规式所构造的DFA状态转移表，另一部分是用来驱动该表的总控程序 `yy1ex()`。当主程序需要从输入字符流中识别一个记号时，只需要调用一次 `yy1ex()` 就可以了。为了使用Lex所生成的词法分析器，我们需要将 `lex.yy.c` 程序用C编译器进行编译，并将相关支持库函数连入目标代码。

LEX源程序必须按照下面的语言规范来写，

1	定义段
2	%%
3	词法规则段
4	%%
5	辅助函数段

注意，辅助函数并非必要的，并且每两个部分之间需要用 `%%` 分隔。Lex源程序中可以有注释，注释由 `/*` 和 `*/` 括起，但是请注意，注释的行首需要有前导空白。

## 定义段

定义段可以分为两部分：

第一部分需要写在 `%{` 和 `%}` 之中，里面以C语法写一些定义与声明。例如，一些头文件，函数声明，宏定义，以及一些全局变量和外部变量定义等等。

本次实验中该部分具体代码为

```
1  %{
2  #include <stdio.h>
3  #include <string.h>
4
5  const int key_num=32;
6  const char* keySet[]={
7      "auto", "break", "case", "char", "const", "continue",
8      "default", "do", "double", "else", "enum", "extern",
9      "float", "for", "goto", "if", "int", "long", "register",
10     "return", "short", "signed", "sizeof", "static", "struct",
11     "switch", "typedef", "union", "unsigned", "void", "volatile",
12     "while"
13 };
14
15 int iskey(char* str);
16
17 %}
```

这里定义了表示一个保留字的数量的变量 `key_num`，一个记录保留字的集合 `keySet`，还定义了一个辅助函数 `iskey`，用以判断遍历到的文本是否为保留字。本次实验中，我们只考虑了32种保留字，并未考虑C11标准中新添加的几种保留字，当然后续可以补充。

第二部分是正规定义和状态定义。正规定义是为了简化后面词法规则而给部分正规式定义的别名。每条正规定义都要顶行写。当然这部分也可以省略。

本次实验中该部分具体代码为，

```

1 letter [A-Za-z]
2 char '['\']*'
3 digit [0-9]
4 preprocess ^#.*
5 string \"[^\"]*\"/>
6 id ({letter}|_)( {letter}|_|{digit})*
7 operater "!"|"%"|"^"|"&"|"*"|"(")|")|"-"|"+"|"="|"=="|"{"|"}"|"["|"]|"|\\"|"|" ":"|" ";"|" "<"|" ">"|" ,"|" "."|" "/"|" "?"|" !=|" "<="|" ">="|" "&&"|" "|" "+"|" "-"|" "*"|" "/"|" "%"|" "<="|" ">="|" "^="|" "&="
8 int {digit}+
9 float ([+-]?){digit}*\\.({digit}+)?((e|E)[+-]?{digit}+)?
10 errorID {digit}({letter}|_|{digit})*({letter}|_){letter}|_|{digit})*
11 %x BLOCK_COMMENT
```

利用正则表达式可以代表一系列的文本。其中 `%x BLOCK_COMMENT` 用来表示注释部分。

## 词法规则段

词法规则段列出的是词法分析器的需要匹配的正规式，以及匹配规则正规式后需要进行的相关动作。

在本次实验中具体代码如下，

```

1 %%
2
3 {preprocess} {
4     FILE *fp=NULL;
5     fp=fopen("token.txt", "a+");
6     fprintf(fp, "<%s,->\n", yytext);
7     fclose(fp);
8 }
9
10 {string} {
11     FILE *fp=NULL;
12     fp=fopen("token.txt", "a+");
13     fprintf(fp, "<$string,%s>\n", yytext);
14     fclose(fp);
15 }
16
17 {char} {
18     FILE *fp=NULL;
19     fp=fopen("token.txt", "a+");
20     fprintf(fp, "<$char,%s>\n", yytext);
21     fclose(fp);
22 }
23
24 {id} {
25     FILE *fp=NULL;
26     fp=fopen("token.txt", "a+");
27     if(!iskey(yytext))
28         fprintf(fp, "<$ID,%s>\n", yytext);
29     else
30         fprintf(fp, "<%s,->\n", yytext);
31     fclose(fp);
32 }
33

```

```

34 {operator} {
35     FILE *fp=NULL;
36     fp=fopen("token.txt", "a+");
37     fprintf(fp, "<$operator,%s>\n", yytext);
38     fclose(fp);
39 }
40
41 {int} {
42     FILE *fp=NULL;
43     fp=fopen("token.txt", "a+");
44     fprintf(fp, "<$int,%s>\n", yytext);
45     fclose(fp);
46 }
47
48 {float} {
49     FILE *fp=NULL;
50     fp=fopen("token.txt", "a+");
51     fprintf(fp, "<$float,%s>\n", yytext);
52     fclose(fp);
53 }
54
55 {errorID} {
56     FILE *fp=NULL;
57     fp=fopen("token.txt", "a+");
58     fprintf(fp, "There is a wrong id: %s.\n", yytext);
59     fclose(fp);
60 }
61
62 "//" .* \n {
63     FILE *fp=NULL;
64     fp=fopen("token.txt", "a+");
65     fprintf(fp, "single line comment.\n");
66     fclose(fp);
67 }
68
69 "/*" {
70     FILE *fp=NULL;
71     fp=fopen("token.txt", "a+");
72     fprintf(fp, "multiple lines comment.\n");
73     fclose(fp);
74     BEGIN(BLOCK_COMMENT);
75 }
76
77 <BLOCK_COMMENT> "*/" {
78     BEGIN(INITIAL);
79 }
80
81 %%

```

每行都是一条规则，该规则的前一部分是正规式，需要顶行首写，后一部分是匹配该正规式后需要进行的动作，这个动作是用C语法来写的，被包裹在{}之内，被Lex翻译器翻译后会被直接拷贝进 `lex.yy.c`。正规式和语义动作之间要有空白隔开。其中用{}扩住的正规式表示正规定义的名字。

## 辅助函数段

辅助函数的写法为C语言写法，辅助函数一般是在词法规则段中用到的函数。这一部分一般会被直接拷贝到 `lex.yy.c` 中。

在本次实验中，我所使用的辅助函数如下

```
1  %%
2
3  int iskey(char* str){
4      for(int i=0;i<key_num;i++){
5          if(strcmp(keySet[i],str)==0){
6              return i+1;
7          }
8      }
9      return 0;
10 }
11
12 int main(){
13     FILE *fp=NULL;
14     fp=fopen("token.txt","w+");
15     fprintf(fp,"Hello, this is a Lexer designed by zhudp3.\n");
16     fclose(fp);
17     yylex();
18     return 0;
19 }
20
```

## 一些变量和函数

`yyin` 和 `yyout`：这是Lex中本身已定义的输入和输出文件指针。这两个变量指明了lex生成的词法分析器从哪里获得输入和输出到哪里。默认：键盘输入，屏幕输出。

`yytext`：指向当前识别的词法单元（词文）的指针

`yylen`：当前词法单元的长度。

`ECHO`：Lex中预定义的宏，可以出现在动作中，相当于 `fprintf(yyout, "%s", yytext)`，即输出当前匹配的词法单元。

`yylex()`：词法分析器驱动程序，用Lex翻译器生成的lex.yy.c内必然含有这个函数。

`yywrap()`：词法分析器遇到文件结尾时会调用yywrap()来决定下一步怎么做：

若 `yywrap()` 返回0，则继续扫描；若返回1，则返回报告文件结尾的0标记。

由于词法分析器总会调用 `yywrap`，因此辅助函数中最好提供 `yywrap`，如果不提供，则在用C编译器编译 `lex.yy.c` 时，需要链接相应的库，库中会给出标准的 `yywrap` 函数（标准函数返回1）。

## 实验结果

首先，对于LEX的使用，我们需要安装对应的环境，输入下面的命令配置环境，安装依赖

```
1 sudo apt install cpp lld flex bison zlib1g-dev \  
2 clang libclang-dev llvm-dev
```

之后，我们编写一个简单的C语言程序，来进行LEX源程序正确性的验证，

```
1 #include <stdio.h>  
2 #define N 10000  
3  
4 int main()  
5 {  
6     int i;  
7     float j,k;  
8     i=10;  
9     j=1.0;  
10    k=-1.2e-10;  
11    char x='y';  
12    char y='x';  
13    //abcdef  
14    char *s="This is a Lexer designed by zhudp3!";  
15    /*Q2  
16    */  
17    return 0;  
18 }
```

进一步的输入下面的命令是编译运行，

```
zhudp3@ubuntu:~/CP/LexAnalysis$ flex clang.lex  
zhudp3@ubuntu:~/CP/LexAnalysis$ gcc -o clang.out lex.yy.c -lfl  
zhudp3@ubuntu:~/CP/LexAnalysis$ ./clang.out <demo.c
```

最终我们得到相应的输出，

```
1 Hello, this is a Lexer desgined by zhudp3.  
2 <#include <stdio.h>,->  
3 <#define N 10000,->  
4 <int,->  
5 <$ID,main>  
6 <$operator,(>  
7 <$operator,)>  
8 <$operator,{>  
9 <int,->  
10 <$ID,i>  
11 <$operator,;>  
12 <float,->  
13 <$ID,j>  
14 <$operator,,>  
15 <$ID,k>  
16 <$operator,;>  
17 <$ID,i>  
18 <$operator,=>  
19 <$int,10>  
20 <$operator,;>  
21 <$ID,j>  
22 <$operator,=>  
23 <$float,1.0>  
24 <$operator,;>
```



```
25 <$ID,k>
26 <$operator,=>
27 <$float,-1.2e-10>
28 <$operator,;>
29 <char,->
30 <$ID,x>
31 <$operator,=>
32 <$char,'y'>
33 <$operator,;>
34 <char,->
35 <$ID,y>
36 <$operator,=>
37 <$char,'x'>
38 <$operator,;>
39 single line comment.
40 <char,->
41 <$operator,*>
42 <$ID,s>
43 <$operator,=>
44 <$string,"This is a Lexer designed by zhudp3!">
45 <$operator,;>
46 multiple lines comment.
47 <return,->
48 <$int,0>
49 <$operator,;>
50 <$operator,}>
```

可以看到所有TOKEN都被正确输出，并且能够准确识别所有类型。特别的，在输出的第27行可以看到该程序可以识别出科学计数法。整型，浮点型，字符型和字符串都可以准确识别。

甚至单行注释和多行注释也可以准确识别，如下面两图

文件(F) 编辑(E) 选择(S) 查看(V) 转到(G) ... demo.c - LexAnalysis [SSH: zdp] - V...

资源管理器 ... C source.c 2 C demo.c clang.lex token.txt

打开的编辑器 C source.c 2 C demo.c clang.lex token.txt

LEXANALYSIS [SSH: ZDP] clang.lex clang.out C demo.c lex.yy.c source.c 2 token.txt

```
C demo.c > main()
1  #include <stdio.h>
2  #define N 10000
3
4  int main()
5  {
6      int i;
7      float j;
8      i=10;
9      j=1.0;
10     i+=j;
11     /*axx
12     */
13     //int 110t;
14     char c='9';
15     char s="999";
16     return 0;
17 }
```

问题 2 输出 调试控制台 终端 端口 bash + - 删除 上 下

axx

```
zhudp3@ubuntu:~/CP/LexAnalysis$ flex clang.lex
zhudp3@ubuntu:~/CP/LexAnalysis$ gcc -o clang.out lex.yy.c -lfl
zhudp3@ubuntu:~/CP/LexAnalysis$ ./clang.out <demo.c
```

> 大纲

SSH: zdp 2 0 0 Live Share 行 13, 列 7 空格: 4 UTF-8 LF C Linux

```
token.txt
14 <$operator,;>
15 <$ID,i>
16 <$operator,=>
17 <$int,10>
18 <$operator,;>
19 <$ID,j>
20 <$operator,=>
21 <$float,1.0>
22 <$operator,;>
23 <$ID,i>
24 <$operator,+=>
25 <$ID,j>
26 <$operator,;>
27 multiple lines comment.
28 single line comment.
29 <char,->
30 <$ID,c>
31 <$operator,=>
32 <$char,'8'>
33 <$operator,;>
34 <char,->
35 <$ID,s>
36 <$operator,=>
```

```
zhudp3@ubuntu:~/CP/LexAnalysis$ flex clang.lex
zhudp3@ubuntu:~/CP/LexAnalysis$ gcc -o clang.out lex.yy.c -lfl
zhudp3@ubuntu:~/CP/LexAnalysis$ ./clang.out <demo.c
axx
```

可以看到两种不同的注释都被分析成功。

除此之外，我还添加了一些报错方式，例如当变量名不符合C语言规范时，该词法分析器会对应的报错，如下面两图，

Visual Studio Code interface showing a C program in `source.c` and its execution output in the terminal.

**File Explorer (资源管理器):**

- 打开的编辑器 (Opened Editors):
  - `source.c` (2 tabs)
  - `demo.c`
  - `clang.lex`
  - `token.txt`
- LEXANALYSIS [SSH: ZDP]:
  - `clang.lex`
  - `clang.out`
  - `demo.c`
  - `lex.yy.c`
  - `source.c` (2 tabs)
  - `token.txt`

**Editor (source.c):**

```
1 #include <stdio.h>
2 #define N 10000
3
4 int main()
5 {
6     int i;
7     float j;
8     i=10;
9     j=1.0;
10    i+=j;
11    /*axx
12    */
13    int 110t;
14    char c='8';
15    char s="999";
16    return 0;
17 }
```

**Terminal:**

```
zhudp3@ubuntu:~/CP/LexAnalysis$ flex clang.lex
zhudp3@ubuntu:~/CP/LexAnalysis$ gcc -o clang.out lex.yy.c -lfl
zhudp3@ubuntu:~/CP/LexAnalysis$ ./clang.out <demo.c
```

**Output:**

```
axx
```

**Status Bar:** SSH: zdp | 2 | 0 | 0 | Live Share | 行 6, 列 11 | 制表符长度: 4 | UTF-8 | LF | C | Linux

```
21 <$float,1.0>
22 <$operator,;>
23 <$ID,i>
24 <$operator,+=>
25 <$ID,j>
26 <$operator,;>
27 multiple lines comment.
28 <int,->
29 There is a wrong id: 110t.
30 <$operator,;>
31 <char,->
32 <$ID,c>
33 <$operator,=>
34 <$char,'8'>
35 <$operator,;>
36 <char,->
37 <$ID,s>
38 <$operator,=>
39 <$string,"999">
40 <$operator,;>
41 <return,->
42 <$int,0>
43 <$operator,;>
```

```
zhudp3@ubuntu:~/CP/LexAnalysis$ flex clang.lex
zhudp3@ubuntu:~/CP/LexAnalysis$ gcc -o clang.out lex.yy.c -lfl
zhudp3@ubuntu:~/CP/LexAnalysis$ ./clang.out <demo.c
```

综上所述，我们可以基本认定实验成功。

## 实验总结

本次实验是编译原理第一次实验，也正是帮助我复习理论中正规式和有限自动机，以及其确定化，最小化的内容。在本次实验中，我选择使用LEX完成项目，这使得我学习到了LEX的相关语法。此外，我还惊叹于LEX的强大计算能力，更加提醒我要提高自己的学习能力和编程能力。

但是，我在本次实验中还存在一些不好的地方。例如，在文件写入和写出的时候，我使用的是 `#include <stdio.h>` 中的文件流，并没有使用使用LEX中自带的 `yyin` 和 `yyout`。这可能引发一些文件安全性的错误。