

基础知识

在本次实验中，我们将用到 Flex 和 Bison 两个工具以及 Cminus-f 语言。这里对其进行简单介绍。

Cminus-f词法

Cminus 是C语言的一个子集，该语言的语法在《编译原理与实践》第九章附录中有详细的介绍。而 Cminus-f 则是在 Cminus 上追加了浮点操作。

1. 关键字

```
1 else if int return void while float
```

2. 专用符号

```
1 + - * / < <= > >= == != = ; , ( ) [ ] { } /* */
```

3. 标识符ID和整数NUM，通过下列正则表达式定义：

```
1 letter = a|...|z|A|...|Z
2 digit = 0|...|9
3 ID = letter+
4 INTEGER = digit+
5 FLOAT = (digit+. | digit*.digit+)
```

4. 注释用 /*...*/ 表示，可以超过一行。注释不能嵌套。

```
1 /*...*/
```

Cminus-f语法

本小节将给出Cminus-f的语法，该语法在Cminus语言（《编译原理与实践》第九章附录）的基础上增加了float类型。

我们将 Cminus-f 的所有规则分为五类。

1. 字面量、关键字、运算符与标识符

- type-specifier
- relop
- addop
- mulop

2. 声明

- declaration-list
- declaration
- var-declaration
- fun-declaration
- local-declarations

3. 语句

- `compound-stmt`
- `statement-list`
- `statement`
- `expression-stmt`
- `iteration-stmt`
- `selection-stmt`
- `return-stmt`

4. 表达式

- `expression`
- `simple-expression`
- `var`
- `additive-expression`
- `term`
- `factor`
- `integer`
- `float`
- `call`

5. 其他

- `params`
- `param-list`
- `param`
- `args`
- `arg-list`

起始符号是 `program`。语法中用到的 token 均以下划线和粗体标出。

1. $\text{program} \rightarrow \text{declaration-list}$
2. $\text{declaration-list} \rightarrow \text{declaration-list } \text{declaration} \mid \text{declaration}$
3. $\text{declaration} \rightarrow \text{var-declaration} \mid \text{fun-declaration}$
4. $\text{var-declaration} \rightarrow \text{type-specifier } \underline{\text{ID}} \text{ ; } \mid \text{type-specifier } \underline{\text{ID}} \text{ [} \underline{\text{INTEGER}} \text{] } ;$
5. $\text{type-specifier} \rightarrow \underline{\text{int}} \mid \underline{\text{float}} \mid \underline{\text{void}}$
6. $\text{fun-declaration} \rightarrow \text{type-specifier } \underline{\text{ID}} \text{ (} \text{params} \text{) } \text{compound-stmt}$
7. $\text{params} \rightarrow \text{param-list} \mid \underline{\text{void}}$
8. $\text{param-list} \rightarrow \text{param-list } \text{ , } \text{param} \mid \text{param}$
9. $\text{param} \rightarrow \text{type-specifier } \underline{\text{ID}} \mid \text{type-specifier } \underline{\text{ID}} \text{ [} \text{]}$
10. $\text{compound-stmt} \rightarrow \{ \text{local-declarations statement-list} \}$
11. $\text{local-declarations} \rightarrow \text{local-declarations } \text{var-declaration} \mid \text{empty}$
12. $\text{statement-list} \rightarrow \text{statement-list } \text{statement} \mid \text{empty}$
 $\text{statement} \rightarrow \text{expression-stmt}$
 $\quad \mid \text{compound-stmt}$
13. $\quad \mid \text{selection-stmt}$
 $\quad \mid \text{iteration-stmt}$
 $\quad \mid \text{return-stmt}$
14. $\text{expression-stmt} \rightarrow \text{expression } ; \mid ;$
 $\text{selection-stmt} \rightarrow \underline{\text{if}} \text{ (} \text{expression} \text{) } \text{statement}$
15. $\quad \mid \underline{\text{if}} \text{ (} \text{expression} \text{) } \text{statement } \underline{\text{else}} \text{statement}$
16. $\text{iteration-stmt} \rightarrow \underline{\text{while}} \text{ (} \text{expression} \text{) } \text{statement}$
17. $\text{return-stmt} \rightarrow \underline{\text{return}} \text{ ; } \mid \underline{\text{return}} \text{expression ;}$

18. $\backslash\text{expression} \rightarrow \text{var} \equiv \text{expression} \mid \text{simple-expression}\backslash$
19. $\backslash\text{var} \rightarrow \underline{\underline{\text{ID}}} \mid \underline{\underline{\text{ID}}} \underline{\quad} [\text{expression}] \backslash$
20. $\backslash\text{simple-expression} \rightarrow \text{additive-expression} \text{ relop } \text{additive-expression} \mid \text{additive-expression}\backslash$
21. $\backslash\text{relop} \rightarrow \underline{\leq} \mid \leq \mid \geq \mid \underline{\geq} \mid \equiv \equiv \mid \underline{\neq} \backslash$
22. $\backslash\text{additive-expression} \rightarrow \text{additive-expression} \text{ addop } \text{term} \mid \text{term}\backslash$
23. $\backslash\text{addop} \rightarrow \underline{+} \mid \underline{-} \backslash$
24. $\backslash\text{term} \rightarrow \text{term} \text{ mulop } \text{factor} \mid \text{factor}\backslash$
25. $\backslash\text{mulop} \rightarrow \underline{*} \mid \underline{/} \backslash$
26. $\backslash\text{factor} \rightarrow (\underline{\quad} \text{expression} \underline{\quad}) \mid \text{var} \mid \text{call} \mid \text{integer} \mid \text{float}\backslash$
27. $\backslash\text{integer} \rightarrow \underline{\underline{\text{INTEGER}}}\backslash$
28. $\backslash\text{float} \rightarrow \underline{\underline{\text{FLOATPOINT}}}\backslash$
29. $\backslash\text{call} \rightarrow \underline{\underline{\text{ID}}} \underline{\quad} (\text{args}) \backslash$
30. $\backslash\text{args} \rightarrow \text{arg-list} \mid \text{empty}\backslash$
31. $\backslash\text{arg-list} \rightarrow \text{arg-list} \text{ , } \text{expression} \mid \text{expression}\backslash$

Flex用法简介

FLEX 是一个生成词法分析器的工具。利用 FLEX，我们只需提供词法的正则表达式，就可自动生成对应的C代码。整个流程如下图：

首先，FLEX从输入文件*.lex或者stdio读取词法扫描器的规范，从而生成C代码源文件lex.yy.c。然后，编译lex.yy.c并与-lfl库链接，以生成可执行的a.out。最后，a.out分析其输入流，将其转换为一系列token。

简答的说，Flex 根据用户定义的正则表达式对输入的字符串进行分析，生成token stream。在我们的编译原理实验中，token stream将被用于后续的语法树生成等后续工作。一个简单的示意图如下：

[illegible]

我们以一个简单的单词数量统计的程序 `wc.1` 为详细介绍下 `Flex` 的功能和用法（请仔细看程序中的注释内容）：

```
1 %option noyywrap
2 %{
3 //在%{和%}中的代码会被原样照抄到生成的lex.yy.c文件的开头，也就是在%{和}%中，你应该按C语言写代码，在这里可以完成变量声明与定义、相关库的导入和函数定义
4 #include <string.h>
5 int chars = 0;
6 int words = 0;
7 %}
8
9 %%
10 /* 注意这里的%%开头*/
11 /* %%开头和%%结尾之间的内容就是使用flex进行解析的部分 */
12 /* 你可以按照我这种方式在这个部分写注释，注意注释最开头的空格，这是必须的 */
```

```

13  /* 你可以在这里使用你熟悉的正则表达式来编写模式，  你可以用C代码来指定模式匹配时对应的动作
    */
14  /* 在%%和%%之间，你应该按照如下的方式写模式和动作 */
15  /* 模式  动作 */
16  /* 其中模式就是正则表达式，动作为模式匹配执行成功后执行相应的动作，这里的动作就是相应的代
    码 */
17  /* 你可以仔细研究下后面的例子 */
18  /* [a-zA-Z]+ 为正则表达式，用于匹配大小写字母 */
19  /* {chars += strlen(yytext);words++;} 则为匹配到大小写字母后，执行的动作（代码），
    这里是完成一个字符累加操作 */
20  /* 这里yytext的类型为 char*，  是一个指向匹配到字符串的指针 */
21  /* yytext是flex自动生成的，在%%和%%之中无需额外定义或者声明 */
22
23  /* 一条 模式 + 动作 */
24  [a-zA-Z]+ { chars += strlen(yytext);words++;}
25
26  /* 另一条 模式 + 动作； . 匹配任意字符，这里匹配非大小写字母的其他字符。这里思考一个问
    题，A既可以被[a-zA-Z]+匹配，也可以被.匹配，在这个程序中为什么A优先被[a-zA-Z]+匹配？如果
    你感兴趣可以去看另一个文档 */
27  . {}
28  /* 对其他所有字符，不做处理，继续执行 */
29  /* 注意这里的%%结尾 */
30  %%
31
32  // flex部分结束，这里可以正常写c代码了
33  int main(int argc, char **argv){
34      // yylex()是flex提供的词法分析例程，调用yylex()即开始执行Flex的词法分析，同样的
    yylex()也是flex自行生成的，无需额外定义和生成，默认输入读取stdin
35      // 如果不清楚什么是stdin，可以自己百度查一下
36      yylex();
37      // 输出 words和chars，这些变量在匹配过程中，被执行相应的动作
38      printf("look, I find %d words of %d chars\n", words, chars);
39      return 0;
40  }

```

使用Flex生成lex.yy.c

```

1  $ flex wc.l
2  $ gcc lex.yy.c
3  $ ./a.out
4  hello world
5  ^D
6  look, I find 2 words of 10 chars

```

注: 在以stdin为输入时，需要按下ctrl+D以退出

至此，你已经成功使用Flex完成了一个简单的分析器！

为了对实验有较好的体验，我建议你好好阅读以下两个关于flex文档：

- [Flex matching](#)
- [Flex regular expressions](#)

Bison用法简介

Bison 是一款解析器生成器 (parser generator)，它可以将 LALR 文法转换成可编译的 C 代码，从而大大减轻程序员手动设计解析器的负担。Bison 是 GNU 对早期 Unix 的 Yacc 工具的一个重新实现，所以文件扩展名为 `.y`。（Yacc 的意思是 Yet Another Compiler Compiler。）

识别一个简单的语言

下面我们以一个简单的语言为例，介绍 Bison 的用法。

每个 Bison 文件由 `%%` 分成三部分。

```
1  %{
2  #include <stdio.h>
3  /* 这里是序曲 */
4  /* 这部分代码会被原样拷贝到生成的 .c 文件的开头 */
5  int yylex(void);
6  void yyerror(const char *s);
7  %}
8
9  /* 这些地方可以输入一些 bison 指令 */
10 /* 比如用 %start 指令指定起始符号，用 %token 定义一个 token */
11 %start reimu
12 %token REIMU
13
14 %%
15 /* 从这里开始，下面是解析规则 */
16 reimu : marisa { /* 这里写与该规则对应的处理代码 */ puts("rule1"); }
17       | REIMU { /* 这里写与该规则对应的处理代码 */ puts("rule2"); }
18       ; /* 规则最后不要忘了用分号结束哦~ */
19
20 /* 这种写法表示 ε -- 空输入 */
21 marisa : { puts("Hello!"); }
22
23 %%
24 /* 这里是尾声 */
25 /* 这部分代码会被原样拷贝到生成的 .c 文件的末尾 */
26
27 int yylex(void)
28 {
29     int c = getchar(); // 从 stdin 获取下一个字符
30     switch (c) {
31         case EOF: return YYEOF;
32         case 'R': return REIMU;
33         default: return YYUNDEF; // 报告 token 未定义，迫使 bison 报错。
34         // 由于 bison 不同版本有不同的定义。如果这里 YYUNDEF 未定义，请尝试 YYUNDEFTOK
35         // 或使用一个随意的整数，如 114514 或 19260817。
36     }
37 }
38
39 void yyerror(const char *s)
40 {
41     fprintf(stderr, "%s\n", s);
42 }
43
44 int main(void)
45 {
```

```

45     yyparse(); // 启动解析
46     return 0;
47 }

```

另外有一些值得注意的点：

1. Bison 传统上将 token 用大写单词表示，将 symbol 用小写字母表示。
2. Bison 能且只能生成解析器源代码（一个 `.c` 文件），并且入口是 `yyparse`，所以为了让程序能跑起来，你需要手动提供 `main` 函数（但不一定要在 `.y` 文件中——你懂“链接”是什么，对吧？）。
3. Bison 不能检测你的 action code 是否正确——它只能检测文法的部分错误，其他代码都是原样粘贴到 `.c` 文件中。
4. Bison 需要你提供一个 `yyllex` 来获取下一个 token。
5. Bison 需要你提供一个 `yyerror` 来提供合适的报错机制。

顺便提一嘴，上面这个 `.y` 是可以工作的——尽管它只能接受两个字符串。把上面这段代码保存为 `reimu.y`，执行如下命令来构建这个程序：

```

1  $ bison reimu.y
2  $ gcc reimu.tab.c
3  $ ./a.out
4  R<-- 不要回车在这里按 Ctrl-D
5  rule2
6  $ ./a.out
7  <-- 不要回车在这里按 Ctrl-D
8  Hello!
9  rule1
10 $ ./a.out
11 blablabla <-- 回车或者 Ctrl-D
12 Hello!
13 rule1      <-- 匹配到了 rule1
14 syntax error <-- 发现了错误

```

于是我们验证了上述代码的确识别了该文法定义的语言 `{ "", "R" }`。

Bison 和 Flex 联动

聪明的你应该发现了，我们这里手写了一个 `yyllex` 函数作为词法分析器。而在上文中我们正好使用 flex 自动生成了一个词法分析器。如何让这两者协同工作呢？特别是，我们需要在这两者之间共享 token 定义和一些数据，难道要手动维护吗？哈哈，当然不用！下面我们用一个四则运算计算器来简单介绍如何让 bison 和 flex 协同工作——重点是如何维护解析器状态、`YYSTYPE` 和头文件的生成。

首先，我们必须明白，整个工作流程中，bison 是占据主导地位的，而 flex 仅仅是一个辅助工具，仅用来生成 `yyllex` 函数。因此，最好先写 `.y` 文件。

```

1  /* calc.y */
2  %{
3  #include <stdio.h>
4      int yyllex(void);
5      void yyerror(const char *s);
6  %}
7
8  %token RET

```

```

9  %token <num> NUMBER
10 %token <op> ADDOP MULOP LPAREN RPAREN
11 %type <num> top line expr term factor
12
13 %start top
14
15 %union {
16     char    op;
17     double  num;
18 }
19
20 %%
21
22 top
23 : top line {}
24 | {}
25
26 line
27 : expr RET
28 {
29     printf(" = %f\n", $1);
30 }
31
32 expr
33 : term
34 {
35     $$ = $1;
36 }
37 | expr ADDOP term
38 {
39     switch ($2) {
40     case '+': $$ = $1 + $3; break;
41     case '-': $$ = $1 - $3; break;
42     }
43 }
44
45 term
46 : factor
47 {
48     $$ = $1;
49 }
50 | term MULOP factor
51 {
52     switch ($2) {
53     case '*': $$ = $1 * $3; break;
54     case '/': $$ = $1 / $3; break; // 想想看，这里会出什么问题？
55     }
56 }
57
58 factor
59 : LPAREN expr RPAREN
60 {
61     $$ = $2;
62 }
63 | NUMBER

```

```

64 {
65     $$ = $1;
66 }
67
68 %%
69
70 void yyerror(const char *s)
71 {
72     fprintf(stderr, "%s\n", s);
73 }

```

```

1  /* calc.l */
2  %option noyywrap
3
4  %{
5  /* 引入 calc.y 定义的 token */
6  #include "calc.tab.h"
7  %}
8
9  %%
10
11 \{ { return LPAREN; }
12 \} { return RPAREN; }
13 "+"|"-" { yylval.op = yytext[0]; return ADDOP; }
14 "*"|"|" { yylval.op = yytext[0]; return MULOP; }
15 "[0-9]+|[0-9]+\.[0-9]*|[0-9]*\.[0-9]+" { yylval.num = atof(yytext); return
NUMBER; }
16 " "|\t { }
17 \r\n|\n|\r { return RET; }
18
19 %%

```

最后，我们补充一个 `driver.c` 来提供 `main` 函数。

```

1  int yyparse();
2
3  int main()
4  {
5      yyparse();
6      return 0;
7  }

```

使用如下命令构建并测试程序：


```

1 $ bison -d calc.y
2     (生成 calc.tab.c 和 calc.tab.h。如果不给出 -d 参数，则不会生成 .h 文件。)
3 $ flex calc.l
4     (生成 lex.yy.c)
5 $ gcc lex.yy.c calc.tab.c driver.c -o calc
6 $ ./calc
7 1+1
8 = 2.000000
9 2*(1+1)
10 = 4.000000
11 2*1+1
12 = 3.000000

```

如果你复制粘贴了上述程序，可能会觉得很神奇，并且有些地方看不懂。下面就详细讲解上面新出现的各种构造。

- `YYSTYPE`: 在 bison 解析过程中，每个 symbol 最终都对应到一个语义值上。或者说，在 parse tree 上，每个节点都对应一个语义值，这个值的类型是 `YYSTYPE`。`YYSTYPE` 的具体内容是由 `%union` 构造指出的。上面的例子中，

```

1 %union {
2     char    op;
3     double  num;
4 }

```

会生成类似这样的代码

```

1 typedef union YYSTYPE {
2     char op;
3     double num;
4 } YYSTYPE;

```

为什么使用 `union` 呢？因为不同节点可能需要不同类型的语义值。比如，上面的例子中，我们希望 `ADDOP` 的值是 `char` 类型，而 `NUMBER` 应该是 `double` 类型的。

- `$$` 和 `$1`, `$2`, `$3`, ...: 现在我们来查看如何从已有的值推出当前节点归约后应有的值。以加法为例：

```

1 term : term ADDOP factor
2     {
3         switch $2 {
4             case '+': $$ = $1 + $3; break;
5             case '-': $$ = $1 - $3; break;
6         }
7     }

```

其实很好理解。当前节点使用 `$$` 代表，而已解析的节点则是从左到右依次编号，称作 `$1`, `$2`, `$3`...

- `%type <>` 和 `%token <>`: 注意，我们上面可没有写 `$1.num` 或者 `$2.op` 哦！那么 bison 是怎么知道应该用 `union` 的哪部分值的呢？其秘诀就在文件一开始的 `%type` 和 `%token` 上。

例如，`term` 应该使用 `num` 部分，那么我们就写

```
1 | %type <num> term
```

这样，以后用 `$` 去取某个值的时候，bison 就能自动生成类似 `stack[i].num` 这样的代码了。

`%token<>` 见下一条。

- `%token`：当我们用 `%token` 声明一个 token 时，这个 token 就会导出到 `.h` 中，可以在 C 代码中直接使用（注意 token 名千万不要和别的东西冲突！），供 flex 使用。`%token <op> ADDOP` 与之类似，但顺便也将 `ADDOP` 传递给 `%type`，这样一行代码相当于两行代码，岂不是很赚。
- `yylval`：这时候我们可以打开 `.h` 文件，看看里面有什么。除了 token 定义，最末尾还有一个 `extern YYSTYPE yyval;`。这个变量我们上面已经使用了，通过这个变量，我们就可以在 lexer 里面设置某个 token 的值。

呼.....说了这么多，现在回头看看上面的代码，应该可以完全看懂了吧！这时候你可能才意识到为什么 flex 生成的分析器入口是 `yylex`，因为这个函数就是 bison 专门让程序员自己填的，作为一种扩展机制。另外，bison（或者说 yacc）生成的变量和函数名通常都带有 `yy` 前缀，希望在这里说还不太晚.....

最后还得提一下，尽管上面所讲已经足够应付很大一部分解析需求了，但是 bison 还有一些高级功能，比如自动处理运算符的优先级和结合性（于是我们就不需要手动把 `expr` 拆成 `factor`, `term` 了）。这部分功能，就留给同学们自己去探索吧！