

# 编译器设计实验报告

2019202463李厚润

## 1. 实验目标

- 完成一个简单的基于PL0语言的编译器
- 拓展功能: **支持for语句**

## 2. PL0语法

```

1  <程序> ::= <分程序>.
2  <分程序> ::= [<常量说明部分>][<变量说明部分>][<过程说明部分>]<语句>
3  <常量说明部分> ::= CONST<常量定义>{,<常量定义>;};
4  <常量定义> ::= <标识符>=<无符号整数>
5  <无符号整数> ::= <数字>{<数字>}
6  <变量说明部分> ::= VAR<标识符>{,<标识符>;};
7  <标识符> ::= <字母>{<字母>|<数字>}
8  <过程说明部分> ::= <过程首部><分程序>;<过程说明部分>;
9  <过程首部> ::= PROCEDURE<标识符>;
10 <语句> ::= <赋值语句>|<复合语句>|<条件语句>|<当型循环语句>|<过程调用语句>|<读语句>|<写语句>|<空>
11 <赋值语句> ::= <标识符>:=<表达式>
12 <复合语句> ::= BEGIN<语句>{;<语句>}END
13 <条件> ::= <表达式><关系运算符><表达式>|ODD<表达式>
14 <条件语句> ::= IF<条件>THEN<语句>
15 <表达式> ::= [+|-]<项>{<加法运算符><项>}
16 <项> ::= <因子>{<乘法运算符><因子>}
17 <因子> ::= <标识符>|<无符号整数>|'(<表达式>')'
18 <加法运算符> ::= +|-
19 <乘法运算符> ::= */
20 <关系运算符> ::= =|<|<=|>|>=
21 <当型循环语句> ::= WHILE<条件>DO<语句>
22 <过程调用语句> ::= CALL<标识符>
23 <读语句> ::= READ'(<标识符>{,<标识符>})'
24 <写语句> ::= WRITE'(<表达式>{,<表达式>})'
25 <字母> ::= a|b|...|x|y|z
26 <数字> ::= 0|1|...|8|9

```

### 3. 文件说明

- lexical\_analysis.l 是词法分析部分，和语法分析类似
- syntax\_analysis.y 是语法/语义分析部分

## 4. 词法分析器设计

- 第一次词法分析实验代码在后面几乎没有任何作用，删除重写，所以就不按照之前的讲了。
- 利用正则表达式来匹配标识符、关键词等等字符：
  - 单字符类型

```

1 digit_ch  [0-9]
2 alpha_ch  [A-Za-z]
3 delimiter_ch  (\.\.\.|\[\\,\\.\\:;\\(\\)\\[\\]\\{\\})
4 blank_ch  (" ")|(\t)|(\r)|(\n)|(\a)|(\b)|(\f)|(\v)|(\0)|
(\)\\)
5 special_ch  [\!\\@#\$%\^&\'\"~\|_ ]
6 operator_ch  [\+\\-\\*\\/\\=\\<\\>]|(>=)|(<=)|(:=)|(!=)

```

◦ 组合类型

```

1 keyword
  if|then|while|do|read|write|call|begin|end|const|var|procedure|odd|in
  teger|float|char|integerarray|for|to|downto|repeat|until
2 identifier  ({alpha_ch}|\\_)( {alpha_ch}|{digit_ch}|\\_)*
3 integer  {digit_ch}+
4 float  {integer}\\. {integer}?([E|e][+|-]?{integer})?
5 character  (\'({digit_ch}|{alpha_ch}|[\\,\\.\\:;\\(\\)\\[\\]\\{\\})|
  {blank_ch}|{special_ch}|[\\+\\-\\*\\/\\=\\<\\>])\\')|({integer})
6 constant {integer}

```

- **匹配串并处理**：如果匹配到了其中的某些特定类别字符串，比如关键词、常数、标识符等等，我们需要进一步细化匹配究竟是匹配了这一大类里的具体哪一个，然后return返回；或者我们需要进一步对yyvalue进行赋值。yyvalue的类型是yystype，默认是int类型，其中对类型的修改见.y语法/语义分析文件。其中每次调用一次yylex函数就会return匹配到的一个特定字符串。其中yylex函数具体应用在后续的词法分析和语义分析都会不断调用yylex。

- return的大写关键词本质是宏，定义在syntax\_analysis.y中。

```

◦ 1 {keyword} {
2     if(!strcmp(yytext, "if")) return IF;
3     if(!strcmp(yytext, "then")) return THEN;
4     if(!strcmp(yytext, "while")) return WHILE;
5     if(!strcmp(yytext, "do")) return DO;
6     if(!strcmp(yytext, "read")) return READ;
7     if(!strcmp(yytext, "write")) return WRITE;
8     if(!strcmp(yytext, "call")) return CALL;
9     if(!strcmp(yytext, "begin")) return BEGIN_;
10    if(!strcmp(yytext, "end")) return END;
11    if(!strcmp(yytext, "const")) return CONST;
12    if(!strcmp(yytext, "var")) return VAR;
13    if(!strcmp(yytext, "procedure")) return PROCEDURE;
14    if(!strcmp(yytext, "odd")) return ODD;
15    if(!strcmp(yytext, "integer")) return INTEGER;
16    if(!strcmp(yytext, "integerarray")) return INTEGERARRAY;
17    if(!strcmp(yytext, "for")) return FOR;
18    if(!strcmp(yytext, "to")) return TO;
19    if(!strcmp(yytext, "downto")) return DOWNT0;
20    if(!strcmp(yytext, "repeat")) return REPEAT;
21    if(!strcmp(yytext, "until")) return UNTIL;
22 }
23
24 {identifier} {
25     yylval.name = strdup(yytext);
26     return ID;
27 }
28

```

```

29 {integer} {
30     int t = atoi(yytext);
31     yylval.valInteger = t;
32     return CONSTANT;
33 }

```

- 对于空白符，移动指针足矣

```

1  [" "\t\r] ;
2  [\n] {yrow++;}
3  . ;

```

- 其中定义了 `%option noyywrap` 表示优化选项，不需要 `wrap` 函数。

## 5. 语法分析器的设计

- 由于语法分析器不支持 `[]` 表示可选项，`{}` 表示重复项等符号，以及存在一些冲突，所以需要我们改写。
- 只需要根据老师词法分析的课件上所给的PL0文法改写即可。
- 比如 `<表达式>{,<表达式>}` 我们可以用 `<表达式列表>` 来代替，然后 `<表达式列表> ::= <表达式>  
<表达式列表> ::= <表达式>, <表达式列表>` 即可表示原来的含义。

### 1. 基础支持

```

1  <表达式> ::= <项列表>
2  <表达式> ::= <加法运算符> <项列表>
3
4  <项列表> ::= <项>
5  <项列表> ::= <项> <加法运算符> <项>
6
7  <无符号整数> ::= CONSTANT
8
9  <项> ::= <因子>
10 <项> ::= <因子> <乘法运算符> <因子>
11
12 <因子> ::= <标识符>
13 <因子> ::= <无符号整数>
14 <因子> ::= '(' <表达式> ')'
15
16 <标识符列表> ::= <标识符>
17 <标识符列表> ::= <标识符>, <标识符列表>
18
19 <表达式列表> ::= <表达式>
20 <表达式列表> ::= <表达式>, <表达式列表>

```

### 2. 程序架构

```

1  <程序> ::= <分程序>.
2
3  <分程序> ::= <常量说明部分> <变量说明部分> <过程说明部分> <语句>
4
5  <常量说明部分> ::= <空>
6  <常量说明部分> ::= CONST <常量定义>;
7
8  <变量说明部分> ::= <空>
9  <变量说明部分> ::= VAR <变量定义>;
10

```

```

11 <过程说明部分> ::= <空>
12 <过程说明部分> ::= <过程首部> <分程序>; <过程说明部分>
13
14 <过程首部> ::= PROCEDURE <标识符>;
15
16 <语句列表> ::= <语句>
17 <语句列表> ::= <语句> ; <语句列表>
18
19 <语句> ::= <赋值语句>
20 <语句> ::= <复合语句>
21 <语句> ::= <条件语句>
22 <语句> ::= <while语句>
23 <语句> ::= <过程调用语句>
24 <语句> ::= <for语句>
25 <语句> ::= <读语句>
26 <语句> ::= <写语句>
27 <语句> ::= <空>
28
29 <条件> ::= <表达式> <关系运算符> <表达式>
30 <条件> ::= ODD<表达式>
31
32 <加法运算符> ::= +
33 <加法运算符> ::= -
34 <乘法运算符> ::= *
35 <乘法运算符> ::= /
36
37 <关系运算符> ::= =
38 <关系运算符> ::= #
39 <关系运算符> ::= <
40 <关系运算符> ::= <=
41 <关系运算符> ::= >
42 <关系运算符> ::= >=

```

### 3. 常量变量定义

- ```

1 <常量定义> ::= <普通常量定义>
2 <常量定义> ::= <普通常量定义>, <常量定义>
3
4 <普通常量定义> ::= <标识符> = <无符号整数>
5
6 <变量定义> ::= <标识符>
7 <变量定义> ::= <标识符>, <变量定义>

```

### 4. 语句

- ```

1 <while语句> ::= WHILE <条件> DO <语句>
2
3 <过程调用语句> ::= CALL <标识符>
4
5 <复合语句> ::= BEGIN <语句列表> END
6
7 <读语句> ::= READ'(' <标识符列表> ')'
8
9 <写语句> ::= WRITE'(' <表达式列表> ')'
10
11 <赋值语句> ::= <标识符>:=<表达式>
12

```

```

13 <条件语句> ::= IF <条件> THEN <语句>
14
15 <for语句> ::= FOR <标识符> := <表达式> TO <表达式> DO <语句>
16 <for语句> ::= FOR <标识符> := <表达式> DOWNTO <表达式> DO <语句>

```

## 5. 终结符:

- 用%token标注的就是终结符，我们需要到词法分析器中加以修改，然后每一次调用yylex函数返回一个这样的token，这就是在词法分析器中提到的宏。

```

1 %token IF THEN WHILE DO READ WRITE CALL BEGIN_ END CONST VAR
  PROCEDURE ODD INTEGER FLOAT CHAR INTEGERARRAY FOR TO DOWNTO REPEAT
  UNTIL
2 %token ID CONSTANT PLUS MINUS MULTIPLY DIVIDE EQ GT LT GE LE ASSIGN
  NEQ
3 %token DOTS COMMA DOT SEMICOLON COLON L0 R0 L1 R1 L2 R2

```

- 对于其他的非终结符，我们只需要一个符号来表示，不需要定义%type（这是后续语义分析需要做的事）。然后定义相关的规则动作就可以了。其中这个工作难度低，重复度高，就是按照给的EBNF范式机械翻译，没啥意思。

- 举个栗子：

- 对于文法 <程序> ::= <分程序> .，我们只需要如下翻译即可。其中有两个儿子节点，分别是 <分程序> 和 .。那么我们通过sechild加入即可。
- 由于这个文法是最终的规约，规约开始符号program，所以我们还需要额外的动作：打印树。其他的规约只需要输出利用到的规约文法即可。

```

1 program : subprogram DOT
2         {
3             $$ = newnode("<程序>");
4             setchild($$, $1);
5             setchild($$, $2);
6             fprintf(opraw, "<程序> -> <分程序>.\n");
7             fprintf(optree, "digraph pic{\n");
8             outtree($$, NULL, 0);
9             fprintf(optree, "}\n");
10        }
11        ;

```

- 具体详情操作见syntax\_analysis.y文件

## 5.2 主函数main说明

```

1 int main(int argc, char **argv){
2     yyin = fopen(argv[1], "r");
3     opraw = fopen(argv[2], "w");
4     optree = fopen(argv[3], "w");
5     yyparse();
6     return 0;
7 }

```

- 接受main的参数，对输入输出进行重定向
- 其中yyparse会默认调用yylex函数。

### 5.3语法树说明

- 语法树的每一个节点是node类型
- 其中基本操作是：

```
1  #define YYSTYPE node*
2  node* newnode(char *str){
3      node* tmp = new node;
4      tmp->text = strdup(str);
5      return tmp;
6  }
7
8  void setchild(node *fa, node *ch){
9      (fa->ch).push_back(ch);
10 }
11
12 void outtree(node *cur, node *fa, int faid){
13     int id = ++num;
14     fprintf(optree, "Node%d[label=\"%s\"]\n", id, cur->text);
15     if(fa != NULL)
16         fprintf(optree, "Node%d -> Node%d\n", faid, id);
17     fprintf(optree, "\n");
18     for(auto i: cur -> ch){
19         outtree(i, cur, id);
20     }
21 }
```

- 我们需要重新定义YYSTYPE为node\*类型，方便我们在词法和语法分析中使用。
- 每一个node节点，包含一个字符串信息信息和所有孩子指针。孩子指针是通过vector存储，方便记录不同数量的孩子。

### 5.4语法树的绘制

- 我们利用outtree来输出语法树的文字形式
- 利用graphviz 的 dot 模块绘制树形图，其中按照先后顺序遍历编号，用->控制连边关系。
- label是标签，用来展示在图形中
- 用 dot optree.out -T png -Gdpi=100 -o optree.png 一条命令即可

### 5.5编译方式指令

```
1  bison -vd syntax_analysis.y
2  flex lexical_analysis.l
3  g++ -std=c++11 -o parse lex.yy.c syntax_analysis.tab.c
4  ./parse test.pl0 opraw.out optree.out
5  dot optree.out -T png -Gdpi=100 -o optree.png
```

- 由于我们使用了vector等c++独有结构以及new函数，所以我们需要g++编译。
- 但是由于其实之前的flex和bison只是对c进行支持，在bison的.y文件中对yylex等外部函数或者变量进行调用的时候，需要添加extern "C"字样，例如：

```
1  extern FILE * yyin;extern int yylineno;extern "C" int yylex();
```

## 6. 语义分析

1. **栈值类型更改**: 类似于语法分析, 我们同样需要更改语法分析栈值的类型, 不同于语法分析 (语法分析只需要按部就班分析语法和输出就可以), 语义分析需要相关的综合属性, 所以我们需要更改这个类型: `#define YYSTYPE SuperNode`。Supernode就是我们更改后的栈值类型, 定义如下:

```
1 struct SuperNode{
2     char* name;
3     int typeH; /*the first tyoe*/
4     int typeS; /*the second type*/
5     int codeAddr; /*the location of the first code of a procedure*/
6     int valInteger; /*the value of variable of constant*/
7     Node *list; /*the nextlist*/
8 };
```

2. **非终结符类型说明**: 类似于语义分析对终结符进行一系列声明, 我们还需要对部分非终结符进行类型说明:

```
1 /*the name of identifier*/
2 %type<name> ID
3 /*the type of expression*/
4 %type<types> plusop multop expression termlist term factor relop
5 /*the value of a caontant variable*/
6 %type<valInteger> CONSTANT
7 /*the nextlist, used for backpatch*/
8 %type<list> statement assignsta compsta statementlist ifsta condition
   whilesta forsta progcallsta readsta writesta
9 /*the anchor, used for backpatch*/
10 %type<codeAddr> checkpoint1
```

3. **代码表示**: 我们对于每一条指令代码, f表示指令, l表示层差, a表示偏移, 通过如下形式表达:

```
1 struct Instructions{
2     char *f;
3     int l;
4     int a;
5 }code[MAX_INSTRUCTIONS];
```

4. **代码生成**: 我们通过如下函数生成代码:

```
1 /*generate the code*/
2 void gen(char *x, int l, int a){
3     code[curIns++] = (Instructions){strdup(x),l,a};
4     /*make the f in lower case*/
5     for(int i = 0; code[curIns - 1].f[i] != 0; i++){
6         if(code[curIns - 1].f[i] <= 'Z')
7             code[curIns - 1].f[i] += 'a' - 'A';
8     }
9 }
```

5. **标识符表示**: 我们定义了结构体Symbol来存放标识符, 其中在symbolTable中存放大量标识符信息。

```

1 struct Symbol{
2     char* name;
3     int typeH; /*the first type*/
4     int types; /*the second type*/
5     int level; /*the level it lays in*/
6     int offset;
7     int size;
8     int codeAddr;
9     int valInteger; /*value*/
10 }symbolTable[MAX_SYMBOL_TABLE_ENTRIES];

```

6. **标识符注册**: 我们有了如上的表示记录, 我们就需要专门的注册函数对新出现的标识符进行注册。

```

1 Symbol *registSymbol(char *name, int typeH, int types, int val){
2     /*only the variable will occupy the stack*/
3     if(typeH == VARIABLE_T)
4         curDataStackSize++;
5     Symbol tmp;
6     tmp.name = strdup(name);
7     tmp.typeH = typeH;
8     tmp.types = types;
9     tmp.level = curLevel;
10    tmp.size = 1;
11    if(typeH == VARIABLE_T){
12        /*if it is the first variable or the local first variable in a
13        procedure*/
14        if(curSymbolTable == 0 || symbolTable[curSymbolTable - 1].level
15        != curLevel || symbolTable[curSymbolTable - 1].typeH != VARIABLE_T)
16            tmp.offset = PROG_OCCUPIED;
17        else
18            tmp.offset = symbolTable[curSymbolTable - 1].offset +
19            symbolTable[curSymbolTable - 1].size;
20    }else if(typeH == CONSTANT_T){
21        tmp.offset = 0;
22        tmp.valInteger = val;
23    }
24    symbolTable[curSymbolTable++] = tmp;
25    return &symbolTable[curSymbolTable - 1];
26 }

```

7. **标识符存取、查找**: 通过lookupSymbol来查找标识符, stoid来存储标识符, lodid来取出标识符

```

1 Symbol *lookupSymbol(char *name){
2     for(int i = curSymbolTable - 1; i >= 0; i--){
3         if(strcmp(name, symbolTable[i].name) == 0)
4             return &symbolTable[i];
5     }
6     return NULL;
7 }
8 Symbol* lodid(char *name){
9     //fprintf(stderr, "%s\n", name);
10    Symbol *tmp = lookupSymbol(name);
11    if(tmp == NULL){
12        fprintf(stderr, "404 %s\n", name);
13        exit(0);
14    }
15    if(tmp -> typeH == VARIABLE_T)

```



```

15     gen("LOD", curLevel - tmp->level, tmp->offset);
16     else{
17         /*If the var is a constant var*/
18         gen("LIT", 0, tmp->valInteger);
19     }
20     return tmp;
21 }
22 void stoid(char *name){
23     Symbol *tmp = lookupSymbol(name);
24     if(tmp == NULL){
25         fprintf(stderr, "404 %s\n", name);
26         exit(0);
27     }gen("STO", curLevel - tmp->level, tmp->offset);
28 }

```

8. backpatch代码回填管理：由于条件表达式和控制流跳转都需要回填技术，所以我们需要维护一些链表，在某一时刻填入跳转目标地址。

```

1
2 /*There will be a cycle*/
3 struct Node* linkList(Node *ptr1, Node *ptr2){
4     if(ptr1 == NULL) return ptr2;
5     if(ptr2 == NULL) return ptr1;
6     Node *ptr1n = ptr1 -> next;
7     Node *ptr2n = ptr2 -> next;
8     ptr1 -> next = ptr2n;
9     ptr2n -> pre = ptr1;
10    ptr2 -> next = ptr1n;
11    ptr1n -> pre = ptr2;
12    return ptr1;
13 }
14
15 struct Node* makeList(int *ptr){
16     Node *res = (Node *)malloc(sizeof(Node));
17     res -> ptr = ptr;
18     res -> pre = res -> next = res;
19     return res;
20 }
21
22 struct Node* insList(Node *list, int *ptr){
23     Node *list2 = makeList(ptr);
24     return linkList(list, list2);
25 }
26
27 void backPatch(Node *list, int x){
28     if(list == NULL) return;
29     Node *tmp = list;
30     do{
31         *(tmp->ptr) = x;
32         tmp = tmp->next;
33     }while(tmp != list);
34 }

```

9. 基于SDT进行语法动作翻译：

- 整个语法规约方法和上次几乎一模一样，稍有一丁点改变，这里就不详细列举了。

- 但是其中会加入许多语义动作，而且不同于语法分析的固定翻译，这次会更加具体且会更难，稍不注意就会有bug。这里的内容很多，我们只需要挑一些具有代表性的介绍一下：

```

1  subprogram : const_cv var_cv {
2      /*Figure out the DataStackSize when compiling*/
3      dataStackSize[curLevel] = curDataStackSize;
4      } proc_cv {
5          if(curLevel > 0)
6              /*the first ins of the current procedure*/
7              funcStack[curLevel - 1] -> codeAddr = curIns;
8          else
9              /*directly enter the master procedure*/
10             code[0].a = curIns;
11             /*Initialize the stack*/
12             gen("INT", 0, dataStackSize[curLevel]);
13         } statement {
14             /*backpatch the list of statement*/
15             backPatch($6, curIns);
16             gen("OPR", 0, 0);
17             curDataStackSize = PROG_OCCUPIED;
18             /*delete the symbols occurred in the current but
19 ending procedure*/
20             deleteSymbol();
21         }
22     };

```

```

1  forsta : /*for to formula*/
2      FOR ID ASSIGN expression{
3          stoid($2);
4      } TO checkpoint1{
5          lodid($2);
6      } expression checkpoint1{
7          /*check boundary*/
8          gen("OPR", 0, 10);
9          gen("JPC", 0, 0);
10         } DO statement {
11             /*the next code of statement is modifying the value
12 of ID*/
13             backPatch($13, curIns);
14             lodid($2);
15             gen("LIT", 0, 1);
16             gen("OPR", 0, 2);
17             stoid($2);
18             /*goto the first ckckpt1 and check the
19 boundary*/
20             gen("JMP", 0, $7);
21             /*the nextlist of forsta equals to the nextlist of
22 falselist of expression*/
23             $$ = makeList(&code[($10)+1].a);
24         }
25         /*for downto formula*/
26         | FOR ID ASSIGN expression{
27             stoid($2);
28             } DOWNTO checkpoint1{
29                 lodid($2);
30             } expression checkpoint1{
31                 gen("OPR", 0, 12);

```

```

29         gen("JPC", 0, 0);
30     } DO statement {
31         backPatch($13, curIns);
32         lodid($2);
33         gen("LIT", 0, 1);
34         gen("OPR", 0, 3);
35         stoid($2);
36         gen("JMP", 0, $7);
37         $$ = makeList(&code[($10)+1].a);
38     }
39 ;

```

#### 10. pcode代码与代码解释器：

- 我们需要一个代码解释器生成c语言，然后用c语言进行编译。interpret只需要不足一些结构和主函数即可。

### 7. 编译运行方法：

- 编译命令

```

1 flex -o LA.cpp lexical_analysis.l
2 yacc -d syntax_analysis.y
3 g++ y.tab.c -o y
4 g++ interpret.cpp -o I

```

- 运行方法：

```

1 ./y test.pl0 test.pcode
2 ./I test.pcode

```

### 8.测试文件说明

- test.pl0是一个有三层嵌套的复杂pl0测试代码
- for.pl0是一个测试for语句的代码

**完结！**

**END！**