

浙江大学



本科实验报告

课程名称 : 编译原理
姓 名 : 杨博淳
学 院 : 计算机科学与技术学院
系 : 计算机科学与技术系
专 业 : 计算机科学与技术
学 号 : 3180101730
指导老师 : 鲁东明

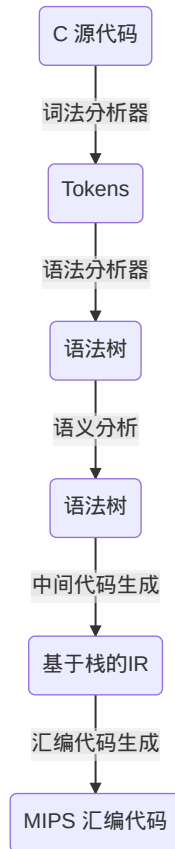
2021年6月19日

编译器

1 总体设计介绍

1.1 流程描述

本项目实现了一个类C语言的编译器，可将C源代码转换为一种栈式机IR并最终编译成MIPS汇编代码。



1.2 语法定义

C文件由外部声明构成，包括函数定义以及全局变量定义，而函数定义包含了函数体，函数体由多个语句组成，语句包括表达式以及控制语句（分支，循环，跳转）。具体语法见语法分析环节。

1.3 运行结果

该项目可在windows或linux下编译(C++ 17)，但最终生成的MIPS汇编需要在MIPS模拟器中运行

使用方法：`mycc src [asm]`

未指定输出文件时，向标准输出打印汇编代码。

指定了输出文件时，向标准输出打印符号表内容，如下：

```
1 | int main() {
2 |     return 0;
3 | }
```

```
(base) → cmake-build-debug git:(main) x ./mygcc ../test/src/test.c ./test.s
symbol table:
name:[main]                type:[Function]                return type:[SInt]                offset:[0]
```

运行MIPS汇编代码或可执行文件：

方法一：在线MIPS模拟器 JsSpim <https://shawnzhong.github.io/JsSpim/>

方法二：qemu-mips

```
1 # 需要docker
2 cd ./test
3 docker pull dockcross/linux-mips:latest
4 docker run --rm dockcross/linux-mips > ./dockercross-mips
5 chmod +x ./dockercross-mips
6 # 使用方法
7 ./dockercross-mips bash -c "\$CC hello.c -o hello; qemu-mips hello"
```

hello world例子：（使用了内嵌函数 `printString`）

```
1 int main() {
2     printString("Hello World!");
3 }
```

汇编代码：

```
(base) → cmake-build-debug git:(main) x ./mygcc ../test/src/test.c
.data
str_0: .asciiz "Hello World!"
.text
.globl main
main:
    addiu $sp, $sp, -4
    sw $ra, 0($sp)
    addiu $sp, $sp, -4
    sw $s8, 0($sp)
    addiu $s8, $sp, 0
    addiu $sp, $s8, 0
    addiu $sp, $sp, -4
    sw $s0, 0($sp)
    la $v0, str_0
    addiu $sp, $sp, -4
    sw $v0, 0($sp)
    lw $a0, 0($sp)
    addiu $sp, $sp, 4
    addiu $sp, $sp, -4
    sw $v0, 0($sp)
    li $v0, 4
    syscall
    lw $v0, 0($sp)
    addiu $sp, $sp, 4
    addiu $sp, $sp, -4
    sw $0, 0($sp)
    addiu $sp, $sp, 4
    addiu $v0, $0, 0
    lw $s0, 0($sp)
    addiu $sp, $sp, 4
    addiu $sp, $s8, 0
    lw $s8, 0($sp)
    addiu $sp, $sp, 4
    lw $ra, 0($sp)
    addiu $sp, $sp, 4
    jr $ra
```

运行结果：

Text Segment

☒ Kernel text
 ☐ Instruction value
 ☒ Source code

User Text Segment

```

[00400000] lw $4, 0($29)           ; 183: lw $a0 0($sp)    # argc
[00400004] addiu $5, $29, 4         ; 184: addiu $a1 $sp 4    # argv
[00400008] addiu $6, $5, 4         ; 185: addiu $a2 $a1 4    # envp
[0040000c] sll $2, $4, 2          ; 186: sll $v0 $a0 2
[00400010] addu $6, $6, $2         ; 187: addu $a2 $a2 $v0
[00400014] jal 0x00400024 [main]   ; 188: jal main
[00400018] nop                   ; 189: nop
[0040001c] ori $2, $0, 10         ; 191: li $v0 10
[00400020] syscall               ; 192: syscall        # syscall 10 (exit)
[00400024] addiu $29, $29, -4      ; 5: addiu $sp, $sp, -4
[00400028] sw $31, 0($29)         ; 6: sw $ra, 0($sp)
[0040002c] addiu $29, $29, -4      ; 7: addiu $sp, $sp, -4
[00400030] sw $30, 0($29)         ; 8: sw $s8, 0($sp)
[00400034] addiu $30, $29, 0       ; 9: addiu $s8, $sp, 0
[00400038] addiu $29, $30, 0       ; 10: addiu $sp, $s8, 0
[0040003c] addiu $29, $29, -4      ; 11: addiu $sp, $sp, -4
[00400040] sw $16, 0($29)         ; 12: sw $s0, 0($sp)
[00400044] lui $2, 4097 [str_0]    ; 13: la $v0, str_0
[00400048] addiu $29, $29, -4      ; 14: addiu $sp, $sp, -4
[0040004c] sw $0, 0($29)         ; 15: sw $0, 0($sp)

```

Execution speed:

Run

Step

Reset

Click line to toggle breakpoint

Output

Hello World!

Log

Based on [SPIM](#) Version 9.1.20 of August 29, 2017 by [James Larus](#).
 Execution finished

1.4 语法特性

本编译器只实现了C的子集，部分功能不支持的。具体支持的语言特性如下。

1.4.1 类型

- 基本类型：int，char，short（有符号）
- 指针
- 数组
- `typedef`
支持定义已实现类型的别名，也可使用 `typedef` 定义的类型来定义复杂类型。
- 全局变量

1.4.2 表达式

- 三元运算表达式：
 - 条件表达式
- 二元运算表达式：

- 加减乘除模
 - 逻辑与、或
 - 按位与、或、异或
 - 左移、右移
 - 关系运算（相等、不等、大小比较）
- 一元运算表达式：
 - 取地址
 - 解引用
 - 正负
 - 按位取反
 - 逻辑非
- 赋值表达式

1.4.3 语句

- 复合语句
 - 包括嵌套作用域带来的同名变量隐藏（variable shadowing）
- 分支：if
- 循环：while , do while , for
- 跳转：continue , break , return

1.4.4 函数

- 多函数
- 递归调用
- 简易的可用于输出字符、字符串、整数的内嵌函数
 - `printChar` , `printString` , `printInt`

2 词法分析

本项目采用了ANTLR (ANother Tool for Language Recognition)作为词法分析和语法分析的工具。

以下定义的词法以及语法都采用了ANTLR的记法。其中词法如下：

```
1  Auto : 'auto';
2  Break : 'break';
3  Case : 'case';
4  Char : 'char';
5  Const : 'const';
6  Continue : 'continue';
7  Default : 'default';
8  Do : 'do';
9  Double : 'double';
10 Else : 'else';
11 Enum : 'enum';
12 Extern : 'extern';
13 Float : 'float';
14 For : 'for';
15 Goto : 'goto';
16 If : 'if';
17 Inline : 'inline';
18 Int : 'int';
19 Long : 'long';
20 Register : 'register';
21 Restrict : 'restrict';
22 Return : 'return';
23 Short : 'short';
24 Signed : 'signed';
25 Sizeof : 'sizeof';
26 Static : 'static';
27 Struct : 'struct';
28 Switch : 'switch';
29 Typedef : 'typedef';
30 Union : 'union';
31 Unsigned : 'unsigned';
32 Void : 'void';
33 Volatile : 'volatile';
34 While : 'while';
35
36 Alignas : '_Alignas';
37 Alignof : '_Alignof';
38 Atomic : '_Atomic';
39 Bool : '_Bool';
40 Complex : '_Complex';
41 Generic : '_Generic';
```

```
42 Imaginary : '_Imaginary';
43 Noreturn : '_Noreturn';
44 StaticAssert : '_Static_assert';
45 ThreadLocal : '_Thread_local';
46
47 LeftParen : '(';
48 RightParen : ')';
49 LeftBracket : '[';
50 RightBracket : ']';
51 LeftBrace : '{';
52 RightBrace : '}';
53
54 Less : '<';
55 LessEqual : '<=';
56 Greater : '>';
57 GreaterEqual : '>=';
58 LeftShift : '<<';
59 RightShift : '>>';
60
61 Plus : '+';
62 PlusPlus : '++';
63 Minus : '-';
64 MinusMinus : '--';
65 Star : '*';
66 Div : '/';
67 Mod : '%';
68
69 And : '&';
70 Or : '|';
71 AndAnd : '&&';
72 OrOr : '||';
73 Caret : '^';
74 Not : '!';
75 Tilde : '~';
76
77 Question : '?';
78 Colon : ':';
79 Semi : ';';
80 Comma : ',';
81
82 Assign : '=';
83 // '*=' | '/=' | '%=' | '+=' | '-=' | '<=>' | '>>=' | '&=' | '^=' | '|='
84 StarAssign : '*=';
85 DivAssign : '/=';
86 ModAssign : '%=';
87 PlusAssign : '+=';
88 MinusAssign : '-=';
89 LeftShiftAssign : '<<=';
```

```

90 RightShiftAssign : '>=';
91 AndAssign : '&=';
92 XorAssign : '^=';
93 OrAssign : '|=';
94
95 Equal : '==';
96 NotEqual : '!=';
97
98 Arrow : '->';
99 Dot : '.';
100 Ellipsis : '...';
101
102 Identifier
103     :   IdentifierNondigit
104         (   IdentifierNondigit
105             |   Digit
106         ) *
107     ;
108
109 fragment
110 IdentifierNondigit
111     :   Nondigit
112         |   UniversalCharacterName
113         |||   // other implementation-defined characters...
114     ;
115
116 fragment
117 Nondigit
118     :   [a-zA-Z_]
119     ;
120
121 fragment
122 Digit
123     :   [0-9]
124     ;
125
126 fragment
127 UniversalCharacterName
128     :   '\\u' HexQuad
129         |   '\\U' HexQuad HexQuad
130     ;
131
132 fragment
133 HexQuad
134     :   HexadecimalDigit HexadecimalDigit HexadecimalDigit HexadecimalDigit
135     ;
136
137 IntegerConstant

```



```
138 : DecimalConstant IntegerSuffix?
139 | OctalConstant IntegerSuffix?
140 | HexadecimalConstant IntegerSuffix?
141 | BinaryConstant
142 ;
143
144 fragment
145 BinaryConstant
146 : '0' [bB] [0-1]+
147 ;
148
149 fragment
150 DecimalConstant
151 : NonzeroDigit Digit*
152 ;
153
154 fragment
155 OctalConstant
156 : '0' OctalDigit*
157 ;
158
159 fragment
160 HexadecimalConstant
161 : HexadecimalPrefix HexadecimalDigit+
162 ;
163
164 fragment
165 HexadecimalPrefix
166 : '0' [xX]
167 ;
168
169 fragment
170 NonzeroDigit
171 : [1-9]
172 ;
173
174 fragment
175 OctalDigit
176 : [0-7]
177 ;
178
179 fragment
180 HexadecimalDigit
181 : [0-9a-fA-F]
182 ;
183
184 fragment
185 IntegerSuffix
```

```

186 : UnsignedSuffix LongSuffix?
187 | UnsignedSuffix LongLongSuffix
188 | LongSuffix UnsignedSuffix?
189 | LongLongSuffix UnsignedSuffix?
190 ;
191
192 fragment
193 UnsignedSuffix
194 : [uU]
195 ;
196
197 fragment
198 LongSuffix
199 : [lL]
200 ;
201
202 fragment
203 LongLongSuffix
204 : 'll' | 'LL'
205 ;
206
207 fragment
208 FloatingConstant
209 : DecimalFloatingConstant
210 | HexadecimalFloatingConstant
211 ;
212
213 fragment
214 DecimalFloatingConstant
215 : FractionalConstant ExponentPart? FloatingSuffix?
216 | DigitSequence ExponentPart FloatingSuffix?
217 ;
218
219 fragment
220 HexadecimalFloatingConstant
221 : HexadecimalPrefix (HexadecimalFractionalConstant | HexadecimalDigitSequence)
  BinaryExponentPart FloatingSuffix?
222 ;
223
224 fragment
225 FractionalConstant
226 : DigitSequence? '.' DigitSequence
227 | DigitSequence '.'
228 ;
229
230 fragment
231 ExponentPart
232 : [eE] Sign? DigitSequence

```

```

233     ;
234
235 fragment
236 Sign
237     :   [+ -]
238     ;
239
240 DigitSequence
241     :   Digit+
242     ;
243
244 fragment
245 HexadecimalFractionalConstant
246     :   HexadecimalDigitSequence? '.' HexadecimalDigitSequence
247     |   HexadecimalDigitSequence '.'
248     ;
249
250 fragment
251 BinaryExponentPart
252     :   [pP] Sign? DigitSequence
253     ;
254
255 fragment
256 HexadecimalDigitSequence
257     :   HexadecimalDigit+
258     ;
259
260 fragment
261 FloatingSuffix
262     :   [f l F L]
263     ;
264
265 fragment
266 CharacterConstant
267     :   '\'' CCharSequence '\''
268     |   'L\'' CCharSequence '\''
269     |   'u\'' CCharSequence '\''
270     |   'U\'' CCharSequence '\''
271     ;
272
273 fragment
274 CCharSequence
275     :   CChar+
276     ;
277
278 fragment
279 CChar
280     :   ~['\\ \r \n]

```

```

281     |   EscapeSequence
282     ;
283 fragment
284 EscapeSequence
285     :   SimpleEscapeSequence
286     |   OctalEscapeSequence
287     |   HexadecimalEscapeSequence
288     |   UniversalCharacterName
289     ;
290 fragment
291 SimpleEscapeSequence
292     :   '\\\' ["?abfnrtv\\]
293     ;
294 fragment
295 OctalEscapeSequence
296     :   '\\\' OctalDigit OctalDigit? OctalDigit?
297     ;
298 fragment
299 HexadecimalEscapeSequence
300     :   '\\\' HexadecimalDigit+
301     ;
302 StringLiteral
303     :   EncodingPrefix? '\'' SCharSequence? '\''
304     ;
305 fragment
306 EncodingPrefix
307     :   'u8'
308     |   'u'
309     |   'U'
310     |   'L'
311     ;
312 fragment
313 SCharSequence
314     :   SChar+
315     ;
316 fragment
317 SChar
318     :   ~["\\\'r\n]
319     |   EscapeSequence
320     |   '\\\'n' // Added line
321     |   '\\\'r\n' // Added line
322     ;
323
324 ComplexDefine
325     :   '#' Whitespace? 'define' ~[#\r\n]*
326     -> skip
327     ;
328

```

```

329 IncludeDirective
330     :   '#' Whitespace? 'include' Whitespace? (('"' ~[\r\n]* '"') | ('<' ~[\r\n]* '>')
    )) Whitespace? Newline
331     -> skip
332     ;
333
334
335 AsmBlock
336     :   'asm' ~{'*' '{' ~'}'* '}'
337     -> skip
338     ;
339
340 LineAfterPreprocessing
341     :   '#line' Whitespace* ~[\r\n]*
342     -> skip
343     ;
344
345 LineDirective
346     :   '#' Whitespace? DecimalConstant Whitespace? StringLiteral ~[\r\n]*
347     -> skip
348     ;
349
350 PragmaDirective
351     :   '#' Whitespace? 'pragma' Whitespace ~[\r\n]*
352     -> skip
353     ;
354
355 Whitespace
356     :   [ \t]+
357     -> skip
358     ;
359
360 Newline
361     :   (   '\r' '\n'?
362         |   '\n'
363         )
364     -> skip
365     ;
366
367 BlockComment
368     :   '/*' .*? '*/'
369     -> skip
370     ;
371
372 LineComment
373     :   '//' ~[\r\n]*
374     -> skip
375     ;

```

3 语法分析

语法如下：

3.1 源文件

```
1 // program file
2 compilationUnit
3     :   translationUnit? EOF
4     ;
5
6 translationUnit
7     :   externalDeclaration+
8     ;
9
10 externalDeclaration
11     :   functionDefinition
12     |   globalDeclaration
13     |   ';' // stray ;
14     ;
15
16 globalDeclaration : declaration;
17
18 functionDefinition
19     :   declarationSpecifiers? declarator '(' parameterTypeList? ')' compoundStatement
20     ;
```

3.2 表达式

按照C运算符的优先级顺序来定义了不同的表达式。

```
1 // expression
2 primaryExpression
3     :   identifier
4     |   constant
5     |   StringLiteral+
6     |   '(' expression ')'
7     ;
8
9 identifier: Identifier;
10
11 constant
12     :   IntegerConstant
13     ;
14
15 postfixExpression
```

```

16      :
17      primaryExpression
18      (
19      '[' expression ']'
20      |
21      '(' argumentExpressionList? ')'
22      )*
23      ;
24
25 argumentExpressionList
26     :   assignmentExpression (',' assignmentExpression)*
27     ;
28
29 unaryExpression
30     :
31     prefixOperator*
32     (
33     postfixExpression
34     |   unaryOperator unaryExpression
35     )
36     ;
37
38 castExpression
39     :   unaryExpression
40     |   '(' typeName pointer* ')' castExpression
41     ;
42
43 multiplicativeExpression
44     :   castExpression (multiplicativeOperator castExpression)*
45     ;
46
47 additiveExpression
48     :   multiplicativeExpression (additiveOperator multiplicativeExpression)*
49     ;
50
51 shiftExpression
52     :   additiveExpression (shiftOperator additiveExpression)*
53     ;
54
55 relationalExpression
56     :   shiftExpression (relationalOperator shiftExpression)*
57     ;
58
59 equalityExpression
60     :   relationalExpression (equalityOperator relationalExpression)*
61     ;
62
63 andExpression

```

```

64     :   equalityExpression (andOperator equalityExpression)*
65     ;
66
67 exclusiveOrExpression
68     :   andExpression (exclusiveOrOperator andExpression)*
69     ;
70
71 inclusiveOrExpression
72     :   exclusiveOrExpression (inclusiveOrOperator exclusiveOrExpression)*
73     ;
74
75 logicalAndExpression
76     :   inclusiveOrExpression (logicalAndOperator inclusiveOrExpression)*
77     ;
78
79 logicalOrExpression
80     :   logicalAndExpression (logicalOrOperator logicalAndExpression)*
81     ;
82
83 conditionalExpression
84     :   logicalOrExpression ('?' expression ':' conditionalExpression)?
85     ;
86
87 assignmentExpression
88     :   unaryExpression assignmentOperator assignmentExpression
89     |   conditionalExpression
90     ;
91
92 expression
93     :   assignmentExpression (',' assignmentExpression)*
94     ;

```

3.3 运算符

```

1 prefixOperator
2     :   '++'
3     |   '--'
4     |   'sizeof'
5     ;
6
7 unaryOperator
8     :   '&' | '*' | '+' | '-' | '~' | '!'
9     ;
10
11 multiplicativeOperator
12     :   '*'

```



```

13     | '/'
14     | '%'
15     ;
16
17 additiveOperator
18     : '+'
19     | '-'
20     ;
21
22 shiftOperator
23     : '<<'
24     | '>>'
25     ;
26
27 relationalOperator
28     : '<'
29     | '>'
30     | '<='
31     | '>='
32     ;
33
34 equalityOperator
35     : '=='
36     | '!='
37     ;
38
39 andOperator : '&';
40
41 exclusiveOrOperator : '^';
42
43 inclusiveOrOperator : '||';
44
45 logicalAndOperator : '&&';
46
47 logicalOrOperator : '|||';
48
49 assignmentOperator : '=' ;

```

3.4 变量声明

```

1 // declaration
2 declaration
3     : declarationSpecifiers initDeclaratorList? ';'
4     ;
5
6 declarationSpecifiers

```

```
7      :   typedefSpecifier? typeName
8      ;
9
10     typeName
11      :   declarationSpecifier+
12      ;
13
14     declarationSpecifier
15      :   typeSpecifier
16      ;
17
18     initDeclaratorList
19      :   initDeclarator (',' initDeclarator)*
20      ;
21
22     initDeclarator
23      :   declarator ('=' initializer)?
24      ;
25
26     typedefSpecifier
27      :   'typedef'
28      ;
29
30     typeSpecifier
31      :   simpleTypeSpecifier
32      ;
33
34     simpleTypeSpecifier
35      :   'void'
36      |   'char'
37      |   'short'
38      |   'int'
39      |   'long'
40      |   'signed'
41      |   'unsigned'
42      ;
43
44     declarator
45      :
46         pointer*
47         directDeclarator
48         array*
49      ;
50
51     array
52      :   ('[' IntegerConstant ']')
53      ;
54
```

```

55 pointer
56     : '*'
57     ;
58
59 directDeclarator
60     : identifier
61     ;
62
63 parameterTypeList
64     : parameterList
65     ;
66
67 parameterList
68     : parameterDeclaration (',' parameterDeclaration)*
69     ;
70
71 parameterDeclaration
72     : typeName declarator
73     ;
74
75 initializer
76     : assignmentExpression
77     ;

```

3.5 语句

```

1 statement
2     : compoundStatement                                #block
3     | expression? ';'                                  #expStmt
4     | 'if' '(' expression ')' statement ('else' statement)? #ifStmt
5     | 'while' '(' expression ')' statement             #whileLoop
6     | 'do' statement 'while' '(' expression ')' ';'    #doWhile
7     | 'for' '(' forCondition ')' statement             #forLoop
8     | 'continue' ';'                                    #continueStmt
9     | 'break' ';'                                       #breakStmt
10    | 'return' expression? ';'                          #returnStmt
11    ;
12
13 compoundStatement
14     : '{' blockItem* '}'
15     ;
16
17 blockItem
18     : statement
19     | declaration
20     ;

```

```

21
22 forCondition
23     :   forInit forCondExpression? ';' forFinalExpression?
24     ;
25
26 forInit:   expression? ';' | declaration;
27 forFinalExpression:   expression;
28 forCondExpression:   expression;

```

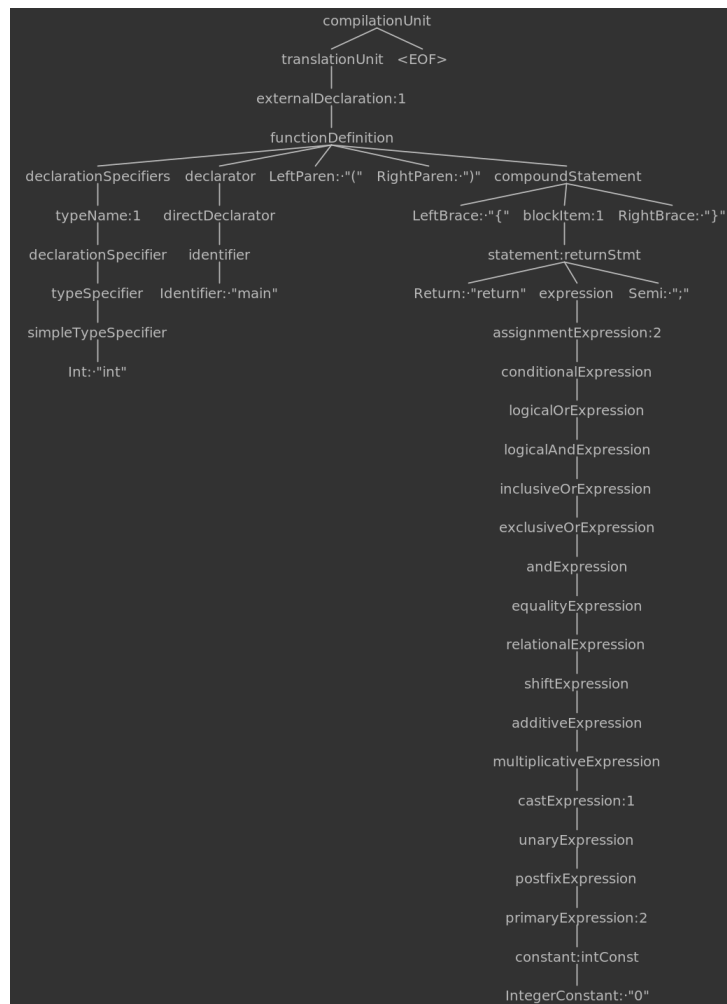
以上定义的词法和语法最终会由ANTLR生成语法树。例如：

```

1  int main() {
2      return 0;
3  }

```

生成的语法树如下：



此外，ANTLR提供了visitor模式的接口，可以很方便地借助visitor模式来访问语法树，进行后续的语义分析和代码生成。

3.6 visitor模式

本项目采用了visitor设计模式。

一般而言，数据结构中保存着元素，需要对元素进行处理。那么为了处理元素，最简单的就是放在数据结构的类中，在类中添加处理的方法。但是如果有很多处理方法，就比较麻烦了，每次增加一种处理方法时可能都需要修改表示数据结构的类。

visitor模式可以解决这个问题，visitor模式将数据结构和处理它的方法分离开。做法是新增一个访问者类，将数据元素的处理交给访问者类，这样以后要新增处理的时候，只需要新增访问者就可以了。

visitor模式很适合用于设计编译器，因为编译过程中的每一趟所对应着对同一个数据结构（语法树）的不同处理。本项目就涉及到两趟，分别执行语义分析以及代码生成，因此也就有两个访问者类：`DeclarationVisitor`和`CodeGenVisitor`。它们继承自基类`CBaseVisitor`，而`CBaseVisitor`是由ANTLR生成的，我们只需要在子类中重载基类的方法就可以访问语法树了。

编译器中的顶层代码如下：输入经过词法分析，语法分析，生成了语法树，然后利用`DeclarationVisitor`与`CodeGenVisitor`访问语法树，分别进行了语义分析与代码生成，最后输出汇编文件。

```
1  CLexer lexer(&input);
2  CommonTokenStream tokens(&lexer);
3  CParser parser(&tokens);
4  tree::ParseTree *tree = parser.compilationUnit();
5  try {
6      DeclarationVisitor declarationVisitor;
7      declarationVisitor.visit(tree);
8      CodeGenVisitor codeGenVisitor;
9      string asm_file = codeGenVisitor.visit(tree).as<string>();
10     std::cout << asm_file << std::endl;
11 }
12 catch (Error &e) {
13     std::cerr << e.what() << std::endl;
14 }
```

4 符号表

符号表项定义：

```
1 class SymTabEntry {
2     public:
3         InitValuePtr initValue;    // init value ctx
4         size_t offset;             // offset in the local var area of stack
5         size_t line;               // line number of this symbol
6         size_t column;             // column number of this symbol
7         string name;
8         CType type;
9     };
```

记录了变量的初始值，栈上的偏移量，定义时的（行，列），变量名，变量类型。

符号表定义：

```
1 class SymTab {
2     unordered_map<string, SymTabEntry> entries;
3     ... (辅助成员变量)
4     class Offsets {
5         size_t words;              // counter
6         size_t offsetInWord;       // counter
7     public:
8         void add(size_t singleSize, size_t count);
9         size_t get() const;
10        void align();
11    };
12    unordered_map<string, Offsets> offsetMap; // func <-> offset
13    ... (方法定义)
14 };
```

首先是一个哈希函数映射，将每个 `id` 映射到一个符号表项。

接下来是与每个函数相关的类 `Offsets`，用于记录每个函数当前栈帧大小的情况，每遇到一个新的变量则需要考虑根据变量大小以及当前栈帧大小来对齐。

5 类型系统

类型是一个树形结构，叶节点为简单类型，如int，char等，父节点则可以是指针、数组、函数等。指针节点的子节点是简单类型表示指向简单类型的指针，而指针、数组类型的子节点同样可以是指针，数组，表示高阶指针或高维数组等。最终的根节点就代表了其类型。

节点基类：

```
1 class CTypeNodeBase {
2 public:
3     virtual BaseType getNodeType() = 0;
4     virtual BaseType getBaseType() = 0;
5     virtual size_t getSize() = 0;
6     virtual CTypeBasePtr typeCast(const CTypeBasePtr &castTo) {
7         return static_pointer_cast<CTypeNodeBase>(castTo);
8     };
9     bool typeCheck(const CTypeBasePtr &srcType);
10    CTypeBasePtr getChild() { return childNode; }
11    size_t getBaseNodeSize();
12 protected:
13    CTypeBasePtr childNode;
14    BaseType nodeType;
15 };
```

每个节点至少有类型（`BaseType`）信息，可能有一个子节点，还需要每个节点的大小，因为数组类型中的每一个维度都需要单独记录一个大小信息。

为了实现`typedef`，还需要记录某一个符号是否是作为`typedef`而定义的。

5.1 基本类型

已实现的基本类型包括：

- signed char
- signed short
- signed int

这类节点只需要知道大小，节点类型以及基节点（子节点）类型都是基本类型。

```
1 class SimpleTypeNode : public CTypeNodeBase {
2 public:
3     BaseType getNodeType() override { return nodeType; }
4     BaseType getBaseType() override { return nodeType; }
5     size_t getSize() override;
6 };
```

5.2 指针

大小为一个字宽（4字节），子节点为其基类型。

```
1 class PointerTypeNode : public CTypeNodeBase {
2 public:
3     BaseType getBaseType() override { return childNode->getBaseType(); }
4     size_t getSize() override { return WORD_BYTES; }
5     BaseType getNodeType() override { return BaseType::Pointer; }
6 };
```

5.3 数组

需要额外的成员变量 `size`，记录数组元素个数。

节点大小为当前节点的 `size` × 子节点的大小。

```
1 class ArrayTypeNode : public CTypeNodeBase {
2     size_t size;
3 public:
4     size_t getSize() override { return size * childNode->getSize(); }
5     size_t getSingleSize() const { return size; }
6     BaseType getNodeType() override { return BaseType::Array; }
7     BaseType getBaseType() override { return childNode->getBaseType(); }
8 };
```

5.4 函数

函数节点的类型记录的是函数的返回类型，另外还需要记录函数的参数列表，函数节点没有大小（在栈帧中不占空间）。

```
1 class FunctionTypeNode : public CTypeNodeBase {
2     vector<CTypeBasePtr> paramList;
3 public:
4     BaseType getNodeType() override { return BaseType::Function; }
5     BaseType getBaseType() override { return childNode->getNodeType(); }
6     size_t getSize() override { return 0; }
7     vector<CTypeBasePtr> &getParamList() { return paramList; }
8 };
```


6 语义分析

在这趟中，编译器将生成符号表，并进行类型检查，对未定义和重复定义的符号报错，以及对不兼容的类型之间的运算进行报错（如不同阶的指针的赋值）。

6.1 标识符

在符号表中查找，如果未定义就报错，返回该标识符的类型。

```
1 // symbol name = function_name@block_order@identifier
2 auto symbol = getCompoundContext() + ctx->Identifier()->getText();
3 // search for symbol
4 auto entry = SymTab::getInstance().get(symbol, ctx->getStart()->getLine(),
5                                     ctx->getStart()->getCharPositionInLine());
6 auto isTypedef = entry.type.isTypedef();
7 return RetType(entry.type.getTypeTree(), !isTypedef, isTypedef);
```

6.2 声明

获取变量的类型标识符，如果包含初始值，则需要检查初始值的类型与变量类型是否匹配，最后将这些信息记录到符号表中。

如果遇到数组声明，则需要检查数组大小是否为常数。

```
1 auto line = ctx->getStart()->getLine();
2 auto column = ctx->getStart()->getCharPositionInLine();
3 auto type = visit(ctx->declarationSpecifiers()).as<CType>();
4 auto declarators = visit(ctx->initDeclaratorList()).as<InitDeclarators>();
5 auto compoundCtx = getCompoundContext();
6 for (auto &declarator : declarators) {
7     auto symbol = compoundCtx + declarator.name;
8     auto typeTree = type.getTypeTree();
9     // array definition
10    for (int i = (int) declarator.arraySizes.size() - 1; i >= 0; i--) {
11        typeTree = static_pointer_cast<CTypeNodeBase>(
12            getArrayType(typeTree, declarator.arraySizes[i]));
13    }
14    // pointer definition
15    for (int i = 0; i < declarator.pointerNum; i++) {
16        typeTree = static_pointer_cast<CTypeNodeBase>(getPointerType(typeTree));
17    }
18    // initial value
19    if (declarator.initValue) {
20        RetType initValueType = visit(declarator.initValue).as<RetType>();
21        // type check
22        if (!typeTree->typeCheck(initValueType.type)) {
23            throw IncompatibleType();
24        }
25    }
26}
```

```

24     }
25 }
26 // add to symbol table
27 SymTab::getInstance().add(symbol, CType(typeTree, bool(type.isTypedef())), line,
column, declarator.initValue, false, typeTree->getNodeType() == BaseType::Array);
28 }
29 return RetType(NoneTypePtr());

```

其中需要对初始值进行类型检查：

```

1 bool CTypeNodeBase::typeCheck(const CTypeBasePtr &srcType) {
2     if (nodeType == BaseType::Pointer && srcType->nodeType == BaseType::Pointer) {
3         return childNode->typeCheck(srcType->childNode);
4     } else if (nodeType == BaseType::Array && srcType->nodeType == BaseType::Array) {
5         if (dynamic_cast<ArrayTypeNode*>(this)->getSingleSize() ==
6             dynamic_pointer_cast<ArrayTypeNode>(srcType)->getSingleSize())
7             return childNode->typeCheck(srcType->childNode);
8     } else if (nodeType == BaseType::Pointer && srcType->nodeType == BaseType::Array) {
9         return childNode->nodeType != BaseType::Pointer;
10    } else if (nodeType == BaseType::Error && srcType->nodeType == BaseType::Error) {
11        return false;
12    } else if (nodeType == srcType->nodeType) {
13        return true;
14    }
15    return false;
16 }

```

6.3 函数定义

需要记录函数的返回类型和参数列表，同样记录到符号表中。

```

1 auto returnType = visit(ctx->declarationSpecifiers()).as<CType>();
2 auto name = visit(ctx->declarator()).as<string>();
3 curFunc = name;
4 // collect parameters information
5 vector<CTypeBasePtr> paramTypes{};
6 if (auto paramTypeList = ctx->parameterTypeList()) {
7     if (auto paramList = paramTypeList->parameterList()) {
8         paramTypes = visit(paramList).as<vector<CTypeBasePtr>>();
9     }
10 }
11 // create a function type, add it to the symbol table
12 auto funcType = CType(static_cast<CTypeBasePtr>(new
FunctionTypeNode(returnType.getTypeTree(), paramTypes)));
13 SymTab::getInstance().add(name, funcType, ctx->getStart()->getLine(), ctx->getStart()-
>getCharPositionInLine());

```

```

14 // visit function body
15 visit(ctx->compoundStatement());
16 blockOrderStack.clear();
17 blockOrder = 0;
18 curFunc = "";
19 return RetType(NoneTypePtr());

```

6.4 复合表达式

遇到复合表达式需要记录当前作用域嵌套的层数，用于进行同名变量的隐藏（shadowing），具体而言是维护一个栈，进入新的作用域时压栈，离开的时候弹出，并且记录到符号表中的变量名携带了作用域的信息。

```

1 blockOrderStack.push_back(blockOrder);
2 blockOrder = 0;
3 auto res = visitChildren(ctx);
4 blockOrder = blockOrderStack.back() + 1;
5 blockOrderStack.pop_back();
6 return res;

```

例如：

```

1 int main() {
2     int x;
3     {
4         x = 1;
5         int x;
6         x = 2;
7     }
8 }

```

编译得到的符号表内容是：

1	name: [main@0@0@x]	type: [SInt]	offset: [1]
2	name: [main@0@x]	type: [SInt]	offset: [0]
3	name: [main]	type: [Function]	return type: [SInt] offset: [0]

`main@0@0@x` 表示 `main` 函数的第 0 个作用域内的第 0 个作用域内的 `x`。

那么在编译的时候，第 2 行定义了 `main@0@x`，然后进入新的作用域，接下来的变量前缀是 `main@0@0@`

到第 4 行时，查找的变量名是 `main@0@0@x`，但是第二个 `x` 还没定义，因此 `main@0@0@x` 不存在，那么将查找外层作用域的变量，即 `main@0@x`，于是找到了第 2 行定义的 `x`。接下来第 5 行定义了新的 `x`，由于没有离开作用域，第 6 行同样会查找 `main@0@0@x`，这时候就会对应到第 5 行定义的 `x`。

6.5 一元运算

一元运算比较特殊的是取地址以及解引用运算符。

取地址运算符只能用在左值上，并且返回的类型是指向其操作数类型的指针。

解引用运算符则只能用在指针类型上，返回的是操作数所指向的类型。

因此就需要得到每个表达式的类型，并且对左值加以区分。

```
1 auto uExp = ctx->unaryExpression();
2 auto op = ctx->unaryOperator();
3 auto rType = visit(uExp).as<RetType>();
4 if (op->Minus() || op->Not() || op->Tilde() || op->Plus()) {
5     return RetType(CTypeBasePtr(rType.type)); // +x, !x, ~x, -x
6 } else if (op->Star()) { // *x, get value of a pointer
7     if (rType.type->getNodeTypes() == BaseType::Pointer) {
8         return RetType(CTypeBasePtr(rType.type->getChild()), true);
9     } else
10        throw InvalidDereference(uExp->getText());
11 } else if (op->And()) { // &x, get address of a left value
12     if (rType.isLeftValue) {
13         auto valueType =
14             static_pointer_cast<CTypeNodeBase>(getPointerType(rType.type));
15         return RetType(valueType);
16     } else
17        throw InvalidLvalue(uExp->getText());
18 }
```

6.6 后缀表达式

6.6.1 下标运算

下标运算适用在数组类型的变量之后，并且对于多维数组，可以进行多个连续的下标运算。

例如，以下操作是合法的：

```
1 int a[3][3][3];
2 int *p = (int *)a[2];
3 int *p = a[2][2];
4 a[1][1][1] = 3;
```

也就是说，当下标运算符的数量与数组维度一致时，返回的是一个左值，而如果运算符的数量少于维度，返回的是一个指针，并且是右值（不能对a[2]赋值）。运算符的数量不能超过数组维度。

另外还需要注意下标运算符内的下标表达式必须是整数类型。

```
1 auto exps = ctx->expression();
2 auto id = ctx->primaryExpression()->identifier();
```

```

3  auto arrayEntry = SymTab::getInstance().get(getCompoundContext() + id->getText());
4  auto resultType = arrayEntry.type.getTypeTree();
5  for (auto exp : exps) {
6      // index should be int
7      if (visit(exp).as<RetType>().type->getNodeType() != BaseType::SInt) {
8          throw InvalidArraySize(exp->getText());
9      }
10     if (resultType->getNodeType() == BaseType::Array ||
11         resultType->getNodeType() == BaseType::Pointer)
12         resultType = resultType->getChild();
13     else { // to many brackets
14         throw InvalidArrayList(ctx->getText());
15     }
16 }
17 // it's left value only when it gets the base type of array
18 if (resultType->getNodeType() != BaseType::Array)
19     return RetType(resultType, true);
20 else
21     return RetType(resultType, false);

```

6.6.2 函数调用

函数调用则只能用在函数名后，括号内的是参数列表，参数列表中的变量数量、类型要与函数定义中的形参列表一致。

```

1  auto argExpList = ctx->argumentExpressionList();
2  if (auto id = ctx->primaryExpression()->identifier()) {
3      auto funcName = id->getText();
4      auto funcEntry = SymTab::getInstance().get(funcName);
5      auto typeTree = funcEntry.type.getTypeTree();
6      if (typeTree->getNodeType() != BaseType::Function) {
7          throw InvalidFuncCall(funcName + " is not callable");
8      }
9      auto funcType = dynamic_pointer_cast<FunctionTypeNode>(typeTree);
10     // argument list
11     if (!argExpList.empty()) {
12         auto exps = argExpList[0]->assignmentExpression();
13         if (exps.size() != funcType->getParamList().size())
14             throw InvalidFuncCall("unmatched parameter number");
15         for (int i = 0; i < exps.size(); i++) {
16             auto lType = funcType->getParamList()[i];
17             RetType rType = visit(exps[i]).as<RetType>();
18             // typecheck between arguments and parameters
19             if (!lType->typeCheck(rType.type)) {
20                 throw IncompatibleType();
21             }

```

```
22     }
23     } else if (!funcType->getParamList().empty())
24         throw InvalidFuncCall("unmatched parameter number");
25     return RetType(funcType->getChild());
26 }
```

6.7 条件表达式

需要保证条件表达式的两个分支表达式的类型相同。

```
1 RetType lType = visit(ctx->expression()).as<RetType>();
2 RetType rType = visit(ctx->conditionalExpression()).as<RetType>();
3 if (!lType.type->typeCheck(rType.type)) { // should be the same
4     throw IncompatibleType();
5 }
6 return RetType(lType.type, true);
```

7 代码生成

在这趟中，编译器将利用语义分析阶段得到的符号表进行代码生成，结果是生成IR代码流，并最终生成汇编代码流。

程序分为代码段（.text）与数据段（.data），数据段中存放全局变量。代码生成的过程中直接向输出流 `_code` 与 `_data` 写入。还需要记录一些和作用域、标签计数相关的信息。

```
1  class CodeGenVisitor : public CBaseVisitor {
2      // code gen
3      ostreamstream _code;
4      ostreamstream _data;
5
6      // string literals
7      size_t stringOrder;
8
9      // compound statement
10     string curFunc;
11     vector<size_t> blockOrderStack;
12     size_t blockOrder;
13     string getCompoundContext();
14
15     // labels
16     size_t labelCount;
17     strings breakStack;
18     strings continueStack;
19
20     ...(code gen visit functions)
21 }
```

7.1 基于栈的IR简介

该编译器采取的IR需要维护一个运算栈，与包括三地址码IR等其他IR不同，这种基于栈的IR没有显式的运算操作数，如 `a = b + c`，在生成的IR中的运算指令中没有a，b，c这些符号，只有代表加法的IR。

具体而言，上述表达式生成的IR大致如下：

```
1  load b;
2  load c;
3  add;
4  store a;
```

也就是说，运算操作的操作数在运算栈中。

此运算栈与运行时栈的关系是，运算栈处于运行时栈的最上方。

类似Java Bytecode，Python Bytecode，本项目采用的IR同属于基于栈的IR。

7.2 MIPS汇编简介

MIPS是RISC（精简指令集）中的一种，具有固定长度的指令格式。MIPS对内存的访问只有`lw(lh/lb)`和`sw(sh/sb)`，对应在寄存器与内存之间进行读取与写入，其余的指令，如算术指令的操作数全部都在寄存器中。

本编译器与MIPS标准的调用约定和变量布局基本一致，但也有区别。具体做法是：

- `$v0` 来传递表达式的中间结果以及函数返回值
- 函数参数的传递全部通过栈进行
- 栈帧上的局部变量按照声明顺序的先后来排列，先声明的变量在栈的高地址位置（偏移更小，离帧指针更近）

7.3 变量标识符

在符号表中查找，得到符号表项。

对于局部变量，需要获取变量在栈上的偏移量（默认按4字节对齐），这是在前一个pass，即语义分析时就得到的常量。然后根据偏移量将变量的实际地址压到栈上，在随后进行变量值的读取或写入就可以根据这个地址进行。

也就是说，在遍历到语法树中的标识符这个节点时，会生成一段将变量的地址压入运算栈的代码。

相关的IR：

指令	参数	含义	IR 栈大小变化
frameAddr	非负整数 k	把当前栈帧中偏移量为 k 的地址压入栈中	增加 1
globalAddr	字符串id	将全局变量id的地址压入栈中	增加 1
load	无参数	弹出栈顶作为地址，加载该地址的内容，将值压入栈中	不变
store	无参数	弹出栈顶作为地址，读取新栈顶作为值，将值写入地址	减少 1
pop	无参数	弹出栈顶，忽略得到的值	减少 1
pushReg	寄存器R	将R的值压入栈中	增加 1
popReg	寄存器R	弹出栈顶的值，写入寄存器R	减少 1

假设a在栈上的偏移量为2，那么标识符a生成的IR为`frameAddr 2;`

对于全局变量b，则生成`globalAddr b;`

MIPS汇编：

指令	汇编
frameAddr k	addiu, \$sp, \$sp, -4; addiu, \$t1, \$sp, -4-4*k; sw, \$t1, 0(\$sp)
globalAddr id la \$t0, id;	addiu \$sp, \$sp, -4; sw \$t0, 0(\$sp);
load	lw \$t1, 0(\$sp); lw \$t1, 0(\$t1); sw \$t1, 0(\$sp);
store	lw \$t1, 4(\$sp); lw \$t2, 0(\$sp); addiu \$sp, \$sp, 4; sw \$t1, 0(\$t2);
pop	addiu \$sp, \$sp, 4
pushReg R	lw R, 0(\$sp); addiu \$sp, \$sp, 4
popReg R	addiu \$sp, \$sp, -4; sw R, 0(\$sp)

代码生成阶段也需要区分表达式类型，但不用像语义分析阶段那么详细，只需要区分左值和右值以及数组类型即可。

原因在于，按照上述的代码生成，如 `a+1`，表达式 `a` 的结果是一个左值，不能直接参与+1的运算，需要在运算前插入一段load代码。而像常量间的运算，如 `1+2`，1和2都是右值，就不需要load。而至于为什么不直接加载所有的标识符的值，则是因为有些时候需要的是标识符的地址，如 `a=1` 进行赋值的时候。

IR生成：

```

1  auto symbol = getCompoundContext() + ctx->getText();
2  auto entry = SymTab::getInstance().get(symbol, ctx->getStart()->getLine(),
3                                     ctx->getStart()->getCharPositionInLine());
4  if (entry.type.isTypedef()) // name for typedef, don't generate code for this id
5      return ExpType();
6  if (entry.name.find_last_of('@') == string::npos) {
7      // global var
8      pushGlobalAddr(entry.name);
9  } else {
10     // local var
11     pushFrameAddr(entry.offset);
12 }
13 if (entry.type.getTypeTree()->getNodeType() == BaseType::Array)
14     return ExpType(ExpType::Type::ARR, 4);
15 else
16     return ExpType(ExpType::Type::LEFT, entry.type.getSize());

```

IR生成MIPS汇编：

```

1  // MIPS imm instruction and memory access instruction
2  inline void iType(const string &op, const string &rs, const string &rt, int imm) {
3      _code << "\t" + op + " $" + rs + ", $" + rt + ", " + to_string(imm) + "\n";
4  }
5

```

```

6  inline void memType(const string &op, const string &rs, const string &rt, int offset) {
7      _code << "\t" + op + " $" + rs + ", " + to_string(offset) + "($" + rt + ")\n";
8  }
9
10 // IR to MIPS
11 inline void pushFrameAddr(size_t k) {
12     iType("addiu", "sp", "sp", -WORD_BYTES);
13     iType("addiu", "t1", "s8", -(int) k);
14     memType("sw", "t1", "sp", 0);
15 }
16
17 inline void pushGlobalAddr(const string &symbol) {
18     la("t0", symbol);
19     pushReg("t0");
20 }
21
22 inline void pushReg(const string &reg) {
23     iType("addiu", "sp", "sp", -4);
24     memType("sw", reg, "sp", 0);
25 }
26
27 inline void popReg(const string &reg) {
28     memType("lw", reg, "sp", 0);
29     iType("addiu", "sp", "sp", 4);
30 }
31
32 inline void pop() { iType("addiu", "sp", "sp", 4); }
33
34 inline void load(size_t size) {
35     memType("lw", "t1", "sp", 0);           // addr
36     memType(loadOp(size), "t1", "t1", 0);  // value
37     memType("sw", "t1", "sp", 0);           // load(to stack)
38 }
39
40 inline void store(size_t size) {
41     memType("lw", "t1", "sp", 4);           // value
42     memType("lw", "t2", "sp", 0);           // addr
43     iType("addiu", "sp", "sp", 4);
44     memType(storeOp(size), "t1", "t2", 0); // store(to addr)
45 }
46
47 // auxiliary function for load/store different size of variables
48 static inline string loadOp(size_t size) {
49     return size == 4 ? "lw" :
50         (size == 2 ? "lh" :
51         (size == 1 ? "lb" :
52         "wrong"));
53 }

```

```

54
55 static inline string storeOp(size_t size) {
56     return size == 4 ? "sw" :
57         (size == 2 ? "sh" :
58             (size == 1 ? "sb" :
59                 "wrong"));
60 }

```

7.4 表达式

7.4.1 二元运算表达式

此处先考虑普通的二元表达式运算，关于指针和数组的运算再另外说明。

二元操作的右操作数在栈顶，左操作数在右操作数下面。

运算的结果存放在栈顶，二元表达式的运算结果都是右值。

IR：

指令	参数	含义	IR 栈大小变化
add	无参数	弹出栈顶两个元素，压入它们的和	减少 1
sub、mul、div、rem	无参数	弹出栈顶两个元素，压入它们的计算结果	减少 1
...（其他二元运算）	无参数	同上	减少 1
const	常数c	加载常数c到栈上	增加 1

1+2会生成IR：`const 1; const 2; add;`

MIPS汇编：

指令	汇编
add	lw \$t0, 0(\$sp); addiu \$sp, \$sp, 4; lw \$t1, 0(\$sp); addiu \$sp, \$sp, 4; add \$t0, \$t0, \$t1; addiu \$sp, \$sp, -4; sw \$t0, 0(\$sp)（解释：即从运算栈中弹出两个操作数到\$t0和\$t1，计算后压栈）
...	（其他二元运算类似）
const c	li \$v0, c; addiu \$sp, \$sp, -4; sw \$v0, 0(\$sp)
logical and	sne \$t0, \$t0, 0; sne \$t1, \$t1, 0; and \$t0, \$t0, \$t1;
logical or	or \$t0, \$t0, \$t1; sne \$t0, \$t0, 0;

实际编写的代码考虑了1+2+3+4这样的连续二元运算可能会频繁地在运算栈上push或pop，针对这种情况进行了一些优化，与这里的汇编代码略有不同，但过程是一样的。具体而言，是将中间结果存在另一个寄存器中，而不是压栈然后紧跟着就弹出，因此节省了2次内存访问开销。

IR生成：

由于二元运算表达式的区别只在与操作符，因此可以有一个统一的模板来生成对应IR。

```
1  template<class T1, typename T2>
2  antlrcpp::Any CodeGenVisitor::genBinaryExpression(vector<T1 *> exps, vector<T2 *> ops)
   {
3      auto exp0Type = visit(exps[0]).template as<ExpType>();
4      if (!ops.empty()) {
5          if (exp0Type.type == ExpType::Type::LEFT) {
6              load(exp0Type.size);
7          }
8          popReg("t0");
9          // save $s0, and use it as accumulator in the following binary expressions
10         pushReg("s0");
11         mov("s0", "t0");
12         for (int i = 0; i < ops.size(); i++) {
13             auto expType = visit(exps[i + 1]).template as<ExpType>();
14             if (expType.type == ExpType::Type::LEFT) {
15                 load(expType.size);
16             }
17             popReg("t0");
18             genBinaryExpressionAsm(dynamic_cast<tree::TerminalNode *>
19                                   (ops[i]->children.front()->getSymbol()->getType()));
20         }
21         mov("v0", "s0");
22         popReg("s0");
23         pushReg("v0");
24         return ExpType(ExpType::Type::RIGHT);
25     }
26     return exp0Type;
27 }
```

此外，当遇到指针运算时，需要做额外的处理，即指针与整数进行运算前先将整数 `x4`。

7.4.2 一元运算表达式

与二元运算类似，一元运算的操作数同样是存放在栈上的。运算时会弹出操作数，计算后将计算结果压入栈中。

因此，对于 `+x`，`-x`，`~x`，`!x` 四个一元运算不做过多介绍，主要的工作在于生成解引用和取地址运算符的代码。

解引用运算符 `*x` 比较简单，只需要将运算栈内的操作数视作是地址，再load一次即可，由于操作的合法性在语义分析阶段已经做了，这里就不需要检查操作数的合法性。

取地址运算符 `&x` 也不难。同样的，语义分析阶段已经确保了 `x` 是一个左值，因此生成 `x` 的代码后，在运行时，栈上存放的是x的地址，因此恰好就是 `&x` 的结果。

然后只需要注意，`*x`的结果是一个左值，其余一元运算的结果都是右值。

```
1  ExpType expType = visit(ctx->unaryExpression()).as<ExpType>();
2  if (op->And()) {      // &x
3      return ExpType(ExpType::Type::RIGHT);
4  }
5  popReg("t0");
6  if (expType.type == ExpType::Type::LEFT) {
7      memType(loadOp(expType.size), "t0", "t0", 0);
8  }
9  if (op->Minus()) {      // -x
10     rType3("sub", "v0", "0", "t0");
11     pushReg("v0");
12     return ExpType(ExpType::Type::RIGHT);
13 } else if (op->Not()) { // !x, 0->1, non-0 ->0
14     iType("sltiu", "v0", "t0", 1);
15     pushReg("v0");
16     return ExpType(ExpType::Type::RIGHT);
17 } else if (op->Tilde()) { // ~x, bit-wise not x
18     rType3("nor", "v0", "t0", "0");
19     rType2("not", "v0", "t0");
20     pushReg("v0");
21     return ExpType(ExpType::Type::RIGHT);
22 } else if (op->Plus()) { // +x, x
23     pushReg("t0");
24     return ExpType(ExpType::Type::RIGHT);
25 } else if (op->Star()) { // *x
26     pushReg("t0");
27     return ExpType(ExpType::Type::LEFT, expType.size);
28 }
```

7.4.3 赋值表达式

语义分析阶段已经确保了等号左边是左值，因此只需要生成等号右边表达式的代码，然后生成等号左边的代码，由于等号左边是左值，因此，在运行时，此时栈顶是等号左边的表达式的地址，栈顶往下紧接着是右边的表达式的值，于是只需要生成`store`就可以了。

赋值表达式返回的都是右值（因此`(a=b)=c`是不合法的，因为`a=b`返回右值，但`a=b=c`是合法的，因为`a=b=c`等价于`a=(b=c)`，没有违背任何左右值的定义）

IR生成：

```

1  ExpType rType = visit(ctx->assignmentExpression()).as<ExpType>();
2  if (rType.type == ExpType::Type::LEFT) {
3      load(rType.size);
4  }
5  ExpType lType = visit(ctx->unaryExpression()).as<ExpType>();
6  if (ctx->assignmentOperator()->Assign()) {
7      store(lType.size);
8  }
9  return ExpType(ExpType::Type::RIGHT, lType.size);

```

7.4.4 后缀表达式

7.4.4.1 下标运算

对于下标运算 `a[x]` 要求 `a` 是指针或数组，`x` 是整数。

假设已经生成了下标运算符前的表达式的代码，那么栈上存放的应该是一个地址，代表数组的基地址 `base`。那么对于每一个 `[x]`，生成表达式 `x` 的代码，然后计算一次地址偏移 `base + x*size`，然后继续下一个 `[x]`。需要注意的是每一维数组的 `size` 不同，例如 `int x[3][6][9]`；那么 `x[1]` 的结果是 `x + 1*6*9*4`，`x[1][2]` 的结果是 `x + 1*6*9*4 + 2*9*4`。

此外和语义分析时一样，需要注意当下标运算符的数量与数组维度一致时，结果是一个左值，而如果运算符的数量少于维度，结果是一个指针，并且是右值。

`a[x][y]` 的IR如下：（假设定义了 `a`：`int a[2][3]`，`int` 大小为4字节）

```

1  (IR for a);
2  (IR for x);
3  const 3*4;
4  mul;
5  add;
6  (IR for x);
7  const 4;
8  mul;
9  add;

```

IR生成：

```

1  auto exps = ctx->expression();
2  auto id = ctx->primaryExpression()->identifier();
3  auto arrayEntry = SymTab::getInstance().get(getCompoundContext() + id->getText());
4  auto resultType = arrayEntry.type.getTypeTree();
5  visit(id);
6  if (resultType->getNode() == BaseType::Pointer) {
7      load(WORD_BYTES);

```

```

8   }
9   for (auto exp : exps) {
10      auto expType = visit(exp).as<ExpType>();
11      if (expType.type == ExpType::Type::LEFT) { // index
12          load(expType.size);
13      }
14      resultType = resultType->getChilid();
15      popReg("t3"); // array index
16      li("t2", (int) resultType->getSize()); // size
17      rType2("mult", "t2", "t3"); // index * size
18      rType1("mflo", "t2");
19      popReg("t3"); // base address
20      rType3("add", "t2", "t3", "t2"); // base address + index * size
21      pushReg("t2"); // address of base[index]
22  }
23  if (resultType->getNodeType() != BaseType::Array)
24      // it's left value only when it gets the base type of array
25      return ExpType(ExpType::Type::LEFT, resultType->getSize());
26  else
27      return ExpType(ExpType::Type::RIGHT, 4);

```

7.4.4.2 函数调用

获取该函数对应的符号表项，然后按逆序，生成参数列表里的表达式的代码，将其结果压入栈中。函数调用结束后，栈顶存放的是函数的返回值。

IR：

指令	参数	含义	IR 栈大小变化
call	函数名f	调用f	增加1，栈上多了返回值
ret	无参数	返回调用者	不变

`x=f(x,y);` 的IR为：（假设x偏移为i）

```

1  (IR for y);
2  (IR for x);
3  call f;
4  frameAddr i;
5  store;

```

此外，为了输出字符串，需要实现系统调用，或者实现 `printf`，后者比较繁琐，还需要实现变长参数列表，因此本编译器直接利用系统调用来实现。编译器内嵌了三个函数，分别可以输出字符，整数或字符串（以 `"\0"` 结尾的字符序列）。调用时参数都在栈上传递。

MIPS汇编：

指令	汇编
call jal f; pop;	(销毁实参);
ret	jr \$ra;

栈帧，函数的prologue以及epilogue，在函数定义中进一步说明。

IR生成：

```

1  auto funcName = id->getText();
2  if (isBuildIn(funcName)) {
3      ... // built-in function
4  } else {
5      auto funcEntry = SymTab::getInstance().get(funcName);
6      size_t offsets = 0;
7      if (!argExpList.empty()) {
8          auto exps = argExpList[0]->assignmentExpression();
9          // push arguments
10         for (int i = (int) exps.size() - 1; i >= 0; --i) {
11             auto expType = visit(exps[i]).as<ExpType>();
12             if (expType.type == ExpType::Type::LEFT) {
13                 load(expType.size);
14             }
15             offsets++;
16         }
17     }
18     call(funcName, offsets);
19     pushReg("v0"); // return value
20     if (funcEntry.type.getTypeTree()->getNodeType() == BaseType::Pointer)
21         return ExpType(ExpType::Type::LEFT, WORD_BYTES);
22     else
23         return ExpType(ExpType::Type::RIGHT, dynamic_pointer_cast<FunctionTypeNode>(
24             funcEntry.type.getTypeTree()->getChild()->getSize());
25 }

```

7.4.5 指针与数组运算

数组的前几维也当作是指针，可以参与一些运算，只有以下的指针算数是合法的：

操作数1	运算符	操作数2	结果
整数	+	指针	指针
指针	+	整数	指针
指针	-	整数	指针
指针	-	指针	整数

对于整数与指针间的运算，需要将整数乘上指针基类型的大小后再相加。

对于指针之间的减法，则需要将相减之后的结果除以指针基类型的大小。

例如，假设有 `int *p;`，`p` 的偏移为 `i`，那么 `p+3` 的IR为：

```
1 | frameAddr i;
2 | load;
3 | push 3;
4 | push 4;
5 | mul;
6 | add;
```

7.5 语句

7.5.1 分支：if-else语句

if-else语句的形式化表达为：`'if' '(' expression ')' statement ('else' statement)?`

分支代码在执行时，首先计算 `expression` 的值，结果存在栈上，然后根据栈上的值是否非0来确定是否跳转到另一分支。而且需要指定跳转目标，这通过打标签来实现。而由于程序中可能出现多个if-else语句，因此跳转的目标要区分，这可以通过维护一个标签计数器来确保每个标签都是独一无二的。

相关IR：

指令	参数	含义	IR 栈大小变化
label	字符串	什么也不做，仅标记一个跳转目的地，用参数字符串标识	不变
beqz	同上	弹出栈顶元素，如果它等于零，那么跳转到参数标识的 label 开始执行	减少 1
br	同上	无条件跳转到参数标识的 label 开始执行	不变

假设标签计数器的值为 `i`。

`if(exp) stmt1; else stmt2;` 的IR如下：

```
1 | (IR for exp);
2 | beqz else_i;
3 | (IR for stmt1);
4 | br end_i;
5 | label "else_i";
6 | (IR for stmt2);
7 | label "end_i" ;
```

`if(exp) stmt1;` 的IR如下：

```

1 (IR for exp);
2 beqz end_i;
3 (IR for stmt1);
4 label "end_i" ;

```

IR生成：

```

1 auto expType = visit(ctx->expression()).as<ExpType>();
2 if (expType.type == ExpType::Type::LEFT) {
3     load(expType.size);
4 }
5 if (ctx->Else()) {
6     // with else statement
7     size_t elseBranch = labelCount++;
8     size_t endBranch = labelCount++;
9     beqz("else_" + to_string(elseBranch));
10    visit(ctx->statement(0));
11    j("end_" + to_string(endBranch));
12    label("else_" + to_string(elseBranch));
13    visit(ctx->statement(1));
14    label("end_" + to_string(endBranch));
15 } else {
16     // without else statement
17     size_t endBranch = labelCount++;
18     beqz("end_" + to_string(endBranch));
19     visit(ctx->statement(0));
20     label("end_" + to_string(endBranch));
21 }

```

7.5.2 条件表达式

条件表达式包含了分支结构，因此在这里介绍。

`exp ? a : b;` 的IR如下：可以发现与if-else语句的IR是完全一样的，但是结果不一样，因为条件表达式的IR的结果是栈上会多出表达式a或表达式b的值，而if-else语句的行为则由其两个分支语句进一步确定。

```

1 (IR for exp);
2 beqz else_i;
3 (IR for a);
4 br end_i;
5 label "else_i";
6 (IR for b);
7 label "end_i";

```

IR生成：

```
1 auto expType = visit(ctx->logicalOrExpression()).as<ExpType>();
2 if (expType.type == ExpType::Type::LEFT) {
3     load(expType.size);
4 }
5 size_t elseBranch = labelCount++;
6 size_t endBranch = labelCount++;
7 beqz("else_" + to_string(elseBranch));
8 expType = visit(ctx->expression()).as<ExpType>();
9 if (expType.type == ExpType::Type::LEFT) {
10     load(expType.size);
11 }
12 j("end_" + to_string(endBranch));
13 label("else_" + to_string(elseBranch));
14 expType = visit(ctx->conditionalExpression()).as<ExpType>();
15 if (expType.type == ExpType::Type::LEFT) {
16     load(expType.size);
17 }
18 label("end_" + to_string(endBranch));
19 return ExpType(ExpType::Type::RIGHT);
```

7.5.3 循环：while语句

```
'while' '(' expression ')' statement
```

每次循环，先检查是否满足循环条件，此外还需要确定 `statement` 内的 `break;` 与 `continue;` 的跳转目标，`break` 跳转到循环语句结束的位置，`continue` 跳转到条件检查的位置。然后需要注意，对于嵌套的循环语句，`break` 或 `continue` 都只是针对外层最近的循环语句而言的，因此需要维护一个栈来记录当前曾的跳转目标。同样的，也需要借助标签计数器来保证跳转目标的唯一性。

`while(exp) stmt;` 的IR：

```
1 label "loop_i";
2 (IR for exp);
3 beqz break_i;
4 (IR for stmt);
5 label "continue_i";
6 br loop_i;
7 label "break_i";
```

IR生成：

```

1  label(loopBegin);
2  // loop condition
3  auto expType = visit(ctx->expression()).as<ExpType>();
4  if (expType.type == ExpType::Type::LEFT) {
5      load(expType.size);
6  }
7  beqz(breakPoint);
8  // loop body
9  visit(ctx->statement());
10 label(continuePoint);
11 j(loopBegin);
12 label(breakPoint);
13 // loop end

```

7.5.4 循环：do-while语句

将statement的位置提前，保证至少执行一次即可

```

1  label "loop_i";
2  (IR for stmt);
3  (IR for exp);
4  beqz break_i;
5  label "continue_i";
6  br loop_i;
7  label "break_i";

```

IR生成：

```

1  label(loopBegin);
2  // body
3  visit(ctx->statement());
4  // cond
5  auto expType = visit(ctx->expression()).as<ExpType>();
6  if (expType.type == ExpType::Type::LEFT) {
7      load(expType.size);
8  }
9  beqz(breakPoint);
10 label(continuePoint);
11 j(loopBegin);
12 label(breakPoint);

```

7.5.5 循环：for语句

```
'for' '(' forInit forCondExpression? ';' forFinalExpression? ')' statement
```

for条件的三个部分都有可能是空，但是不影响整体IR的顺序。此外还需要注意，for循环自带了一层作用域，for条件中的声明部分的id是可以与for循环之外的变量id相同的。

```
1 (IR for exp/decl);
2 label "loop_i";
3 (IR for forCondExp);
4 beqz break_i;
5 (IR for stmt);
6 label "continue_i";
7 (IR for forFinal);
8 br loop_i;
9 label "break_i";
```

IR生成：

```
1 if (auto exp = ctx->forCondition()->forInit()->expression()) {
2     visit(exp);
3     pop(); // pop unused value in stack
4 } else if (auto dec = ctx->forCondition()->forInit()->declaration()) {
5     visit(dec);
6 }
7 label(loopBegin);
8 // cond
9 if (auto cond = ctx->forCondition()->forCondExpression()) {
10     auto expType = visit(cond->expression()).as<ExpType>();
11     if (expType.type == ExpType::Type::LEFT) {
12         load(expType.size);
13     }
14     beqz(breakPoint);
15 }
16 // body
17 visit(ctx->statement());
18 label(continuePoint);
19 if (auto final = ctx->forCondition()->forFinalExpression()) {
20     visit(final);
21     pop(); // pop unused value in stack
22 }
23 j(loopBegin);
24 label(breakPoint);
```

7.5.6 跳转：break语句/continue语句

这两个都很简单，只需要指定跳转目标即可，而跳转目标已经在生成循环IR的过程中压入一个栈中了，这时候只需要直接使用，如果栈为空则报错。

```
1 antlrcpp::Any CodeGenVisitor::visitContinueStmt(CParser::ContinueStmtContext *ctx) {
```

```

2     if (continueStack.empty()) {
3         throw InvalidContinue();
4     } else
5         j(continueStack.back());
6     return ExpType();
7 }
8
9 antlrcpp::Any CodeGenVisitor::visitBreakStmt(CParser::BreakStmtContext *ctx) {
10     if (continueStack.empty()) {
11         throw InvalidBreak();
12     } else
13         j(breakStack.back());
14     return ExpType();
15 }

```

7.5.7 复合语句

与语义分析的过程一样，需要维护一个栈来记录当前作用域的信息，每遇到一个复合语句就需要压栈，复合语句结束需要弹出，以便于在复合语句内部查找变量时能正确关联到符号表里的符号。

然后生成复合语句的IR，只需要对复合语句内的语句逐条进行代码生成即可。

IR生成：

```

1 blockOrderStack.push_back(blockOrder);
2 blockOrder = 0;
3 for (auto item : ctx->blockItem()) {
4     visit(item);
5 }
6 blockOrder = blockOrderStack.back() + 1;
7 blockOrderStack.pop_back();

```

7.6 函数定义

7.6.1 栈帧（运行环境）

假设有以下代码：

```

1  int g(int a, int b, int c){
2      ...
3  }
4
5  int f(int a, int b){
6      return 1 + g(x, y, z);
7  }
8
9  int main(){
10     return f(3, 4);
11 }

```

在调用 `g` 时的栈布局如下：



也就是说，每个栈帧顶端是运算栈，然后是局部变量，然后是返回地址，帧指针，然后是函数参数。

7.6.2 prologue

每次调用函数时，调用者首先需要保存返回地址以及之前的帧指针。

然后就需要将栈指针赋值给新的帧指针。

接下来需要处理参数传递。在调用函数之前，需要准备好函数的参数，参数按逆序直接压到栈上，于是在调用时，就可以将参数的值赋值给被调用函数的形参。

紧接着需要在栈上为所有的局部变量申请空间，需要注意，函数形参也是局部变量。

函数的prologue结束了，接下来就是遍历函数体内的每条语句。

IR生成：

```
1  curFunc = ctx->declarator()->directDeclarator()->identifier()->getText();
2  // prologue: allocate an frame
3  label(curFunc);
4  pushReg("ra");
5  pushReg("s8");
6  mov("s8", "sp");
7
8  // load arguments to the new frame
9  auto paramNames = SymTab::getInstance().getParamNames(curFunc);
10 size_t argOffset = 1;
11 for (auto &paramName : paramNames) {
12     auto entry = SymTab::getInstance().get(paramName);
13     auto size = entry.type.getTypeTree()->getSize();
14     argOffset++;
15     auto paramOffset = -(int) entry.offset;
16     // arg -> param
17     memType("lw", "t0", "s8", WORD_BYTES * (int) argOffset);
18     memType(storeOp(size), "t0", "s8", paramOffset);
19 }
20
21 // allocate local vars on stack
22 size_t offsets = SymTab::getInstance().getTotalOffset(curFunc);
23 iType("addiu", "sp", "s8", -(int) offsets);
24
25 // visit function body
26 pushReg("s0");
27 visit(ctx->compoundStatement());
28
29 // exit a function scope
30 blockOrderStack.clear();
31 blockOrder = 0;
32 curFunc = "";
```

7.6.3 epilogue：return语句

函数的epilogue在return时执行，具体的步骤为：

加载返回值，置于栈顶。

销毁当前函数的栈帧，回收空间。

恢复保存的寄存器，包括帧指针和返回地址。

返回调用者。

IR生成：

```
1  // epilogue
2  auto expType = visit(ctx->expression()).as<ExpType>();
3  if (expType.type == ExpType::Type::LEFT) {
4      load(expType.size);
5  }
6  popReg("v0"); // return value
7
8  // restore saved registers
9  popReg("s0");
10
11 // epilogue: deallocate an frame
12 size_t offsets = SymTab::getInstance().getTotalOffset(curFunc);
13
14 // restore old frame pointer and return address
15 mov("sp", "s8");
16 popReg("s8");
17 popReg("ra");
18 ret();
```

7.7 全局变量

全局变量不是放在栈上的，而是放在 `.data` 段。因此，就不能通过计算相对于帧起始地址的方式来获取变量地址。利用之前定义的IR中的 `globalAddr` 就可以实现。

此外，还需要在 `.data` 段中定义全局变量，汇编如下：

```
1  .globl id
2      .align 4
3  id:
4      .space 4
```

`.globl id` 是为了使得 `id` 全局可见，`.align 4` 是为了对齐，标签 `id` 定义了一段大小为4字节的空间。

8 测试

测试脚本见 `test/test.sh`，通过检查 `main` 函数的返回值，也即程序的返回值，将其与 `gcc` 编译的程序的运行结果作对比来测试结果的正确性。

9 总结

总体来说，工作量很大，收获也很大。因为这个项目是我自己完成的，也没有采取LLVM这样的工具来做代码生成，因此可以说，对于编译器的从前端到后段的整个工作流程，我都有了更深刻的理解。只是这个项目距离一个真正完整的C编译器而言还有一段距离，因为包括 `struct`，`enum` 等自定义类型，以及变长参数列表在内的许多语法特性都没有实现。而没有变长参数列表，甚至连 `printf` 函数都没有办法实现。

但总之，通过这次实验，我掌握了编译器设计的基本原理和步骤，包括如何设计一个语言的语法，利用词法和语法分析工具生成一个语言的前端，以及如何实现编译器的后端：即如何进行语义分析，生成符号表，进行类型检查，以如何进行代码生成，如何设计与生成中间表示，以及最后如何将中间表示翻译为汇编语言。还包括了如何设计函数的调用约定，设置运行环境。其间还要自己设计测试代码来验证编译器的正确性等等。