

实验5 目标代码生成

202011081001 卢皓宇 2022年1月5日

实验要求

在词法分析、语法分析、语义分析和中间代码生成程序的基础上，将C++源代码翻译为MIPS32指令序列（可以包含伪指令），并在 SPIM Simulator上运行

面向样例，不好说能不能实际运行/doge

基本要求

将中间代码正确转化为MIPS32汇编代码

“正确”是指该汇编代码在SPIM Simulator（命令行或Qt版本均可）上运行结果正确

附加要求

1. 对寄存器分配进行优化，如使用局部寄存器分配办法等
2. 对指令进行优化

化简要求

- 寄存器的使用与指派可以不必遵循MIPS32的约定。只要不影响在SPIM Simulator中的正常运行，你可以随意分配 MIPS体系结构中的32个通用寄存器，而不必在意哪些寄存器应该存放参数、哪些存放返回值、哪些由调用者负责保存、哪些由被调用者负责保存，等等。
- 栈的管理（包括栈帧中的内容及存放顺序）也不必遵循MIPS32的约定。你甚至可以使用栈以外的方式对过程调用间各种数据的传递进行管理，前提是你输出的目标代码“正确”。

实验分工

- 卢皓宇：基于树形IR完成大多数生成目标代码的程序（模式 -> 代码），并基于样例1和2进行完善和优化；设计查询变量表VarReg以确定变量所属寄存器、暂存当前生成的函数/过程结点以便查询的方法并基本实现；GitHub项目管理
- 姜潮涌：完成确定变量所属寄存器方法findMark的封装，以及RELOP中涉及的常量的管理方法findNum2的封装；完善代码细节
- 段欣然：提出BFS生成各函数/过程的理念并通过队列（链表）实现；完善代码细节
- 杜隆清：修改、增加样例以测试代码鲁棒性，提出修改建议

实验环境

- Bison: v3.5.1
- Flex: v2.6.4
- 虚拟机程序irsim
- 本地程序编写与测试
 - Windows Subsystem for Linux 2: Ubuntu 20.04 LTS
 - gcc: 9.4.0
- 云主机测试

- o BNUCLOUD云主机: Ubuntu 18.04 LTS
- o gcc: 7.5.0

实验设计

指令选择

由于我们在实验四中采用了树形IR，且并未采用其他数据结构对中间代码进行存储，而是直接依托于语法树完成，所以实验五的实现也基本依托于语法树。

相应地，指令选择机制也就是DFS语法树上的结点，以确定该结点及其兄弟孩子是否符合相应的模式。

我们举一个比较复杂的例子：赋值语句，对应到产生式应该是Exp: Exp ASSIGNOP Exp。赋值语句对应指令的可能性相对较多：我们假设已声明变量a，其对应的寄存器是\$a

- a=1 -> li \$a, 1
- a=b -> move \$a, \$b
- a=b+c -> add \$a, \$b, \$c (这个可能的写法比较多)

所以我们需要对ASSIGNOP构建出不同模式，每对应一个模式就生成其相应的代码。下面是我们代码的部分展示

```
1  if(head->child->bro) {
2      if(!strcmp(head->child->bro->name, "ASSIGNOP")) {
3          // Exp ASSIGNOP Exp
4          if(!strcmp(head->child->bro->bro->child->name, "INT")) {
5              // Exp(ID) = Exp(INT)
6              int mark = findMark(head->child->child->id);
7              fprintf(f, "\tli $t%d, %d\n", mark, head->child->bro->bro->child->intValue);
8
9              sprintf(expReg, "t%d", mark);
10             return expReg;
11         }
12         else if(!strcmp(head->child->bro->bro->child->name, "ID") &&
!head->child->bro->bro->child->bro) {
13             // Exp(ID) = Exp(ID)
14             int mark = findMark(head->child->child->id);
15             fprintf(f, "\tmove $t%d, $t%d\n", mark, findMark(head->child->bro->bro->child->id));
16
17             sprintf(expReg, "t%d", mark);
18             return expReg;
19         }
20         else if(head->child->bro->bro->child->bro->bro->bro) { // 一定要
注意先写长的再写短的
21             // ID LP Args RP
22             pushQueue(hed, head->child->bro->bro->child->id);
23             genExp(head->child->bro->bro, f);
24             int mark = findMark(head->child->child->id);
25             fprintf(f, "\tmove $t%d, $v0\n", mark); // 赋值
26
27             sprintf(expReg, "t%d", mark);
28             return expReg;
29         }
```

```

30         else if(head->child->bro->bro->child->bro) {
31             // 运算式/函数调用赋值给变量
32             char tmp[5];
33             strcpy(tmp, genExp(head->child->bro->bro, f));
34
35             sprintf(expReg, "t%d", findMark(head->child->child->id));
36             if(strcmp(expReg, tmp)) { // 不是自增
37                 fprintf(f, "\tmove %s, %s\n", expReg, tmp);
38             }
39             return expReg;
40         }
41         else {
42             printf("no\n");
43         }
44     }

```

其他的目标代码生成也遵照了类似的设计方式。

感觉我们在上一步走偏了，中间代码在这种指令选择机制下好像起不到太大的作用，例如ARG、PARAM、CALL之类的就完全没用到。虽然直接用语法树多少也能做，但这样总感觉差了点意思，而且也不够方便。。。sad

寄存器分配

类似朴素寄存器分配算法。我们没有引入释放寄存器的操作，不过在一些中间变量的处理上（例如上文提到的自增、相同的常量），我们尽可能地节约了寄存器的使用。所以我们也简单介绍一下我们的做法

- 对于变量，我们设置了一个二维char数组VarReg作为变量名表，在建立语法树时就将所有**声明的变量**（函数形参列表除外）按顺序加入该表，以其位置作为对应的寄存器编号；同时引入变量计数Regcnt并同时维护，以确定VarReg的终止位置。当变量被使用时，就从起始位置开始直到Regcnt查找该变量名，并返回其位置作为对应的寄存器编号
- 在前面VarReg的基础上，我们将**第一次被使用的常量**加入该表，但此时不维护Regcnt（为区分变量与常量），后续使用同一常量直接从Regcnt位置开始向后搜索即可找到该常量对应的寄存器编号。当然其实应该也只有RELOP因为伪指令格式要求需要传入常量所属的寄存器

```

1  char expReg[10]; // Exp中间变量的寄存器名
2
3  int findMark(char* id) {
4      int mark = 0;
5      for(int i = 0; i < Regcnt; i++) {
6          if(!strcmp(VarReg[i], id)) {
7              mark = i;
8              break;
9          }
10     }
11     return mark;
12 }
13
14 int findNum2(int n, FILE* f) {
15     if(n == 0){
16         fprintf(f, "\tli $t0, 0\n");
17         return 0;
18     }
19     for(int i = Regcnt; i < 20; i++) {
20         if(ImmReg[i] == n) {

```

```

21         fprintf(f, "\tli $t%d, %d\n", i, n);
22         return i;
23     }
24     else if(ImmReg[i] == 0) {
25         fprintf(f, "\tli $t%d, %d\n", i, n);
26         ImmReg[i] = n;
27         return i;
28     }
29 }
30 }

```

也是感觉这个步骤树形IR不太好操作，毕竟没有将代码设置成一个线性结构，不知道怎么与课上的基本块、活性分析对应起来。。。而且时间和精力确实比较紧张

函数调用

宽搜：确定各函数/过程翻译顺序

基本思路

- 整个生成目标代码是一个BFS函数/过程队列的过程，初始从main开始
- 遇到调用函数（ID LP RP或ID LP Args RP），将该函数push进队列中（pushQueue方法），等待当前过程生成完毕再去生成被调用函数的目标代码
- write/read函数一开始在文件头部生成，故不考虑

```

1  struct Queue {
2      char name[10];
3      struct Queue* nxt;
4  };
5
6  struct Queue* hed; // 队首存main
7
8  void pushQueue(struct Queue* Qhed, char* name) {
9      if(!strcmp(name, "read") || !strcmp(name, "write")) return;
10     struct Queue* tmp;
11     tmp = (struct Queue*)malloc(sizeof(struct Queue*));
12     strcpy(tmp->name, name);
13     tmp->nxt = NULL;
14     struct Queue* q = (struct Queue*)malloc(sizeof(struct Queue*));
15     q = Qhed;
16     while(q->nxt){
17         if(!strcmp(q->name, name)) return;
18         q = q->nxt;
19     }
20     q->nxt = tmp;
21 }
22
23 void bfsGenExtDefList(struct node* head, FILE* fp) {
24     struct Queue* q = (struct Queue*)malloc(sizeof(struct Queue*));
25     currentFunc = (struct node*)malloc(sizeof(struct node));
26     q = hed;
27     struct node p = *head;
28     while(q){
29         strcpy(arg, "");
30         // strcpy(stackRet, "");

```

```

31     // strcpy(regName, "");
32
33     p = *head;
34     while(p.child && strcmp(q->name, p.child->child->bro->child->id)) {
35         p = *(p.child->bro);
36     }
37     currentFunc = p.child->child->bro;
38     genExtDef(p.child, fp);
39     // if(q->nxt)
40     q = q->nxt;
41     // else break;
42 }
43 }
44
45 void genAsm(struct node* head, FILE* fp) {
46     // ...
47
48     hed = (struct Queue*)malloc(sizeof(struct Queue*));
49     strcpy(hed->name, "main");
50     hed->nxt = NULL;
51     bfsGenExtDefList(head->child, fp);
52     free(hed);
53 }
54
55 char* genExp(struct node* head, FILE* fp) {
56     // ...
57     else if(head->child->bro->bro->child->bro->bro->bro) { // 一定要注意先写长的
        的再写短的
58         // ID LP Args RP
59         pushQueue(hed, head->child->bro->bro->child->id);
60         genExp(head->child->bro->bro, f);
61         int mark = findMark(head->child->child->id);
62         fprintf(f, "\tmove $t%d, $v0\n", mark); // 赋值
63
64         sprintf(expReg, "t%d", mark);
65         return expReg;
66     }
67     // ...
68     else if(!strcmp(head->child->bro->name, "LP") && !strcmp(head->child-
        >bro->bro->name, "RP")) {
69         // ID LP RP
70         pushQueue(hed, head->child->id);
71         fprintf(f, "\taddi $sp, $sp, -4\n");
72         fprintf(f, "\tsw $ra, 0($sp)\n");
73         fprintf(f, "\tjal %s\n", head->child->id);
74         fprintf(f, "\tlw $ra, 0($sp)\n");
75         fprintf(f, "\taddi $sp, $sp, 4\n");
76
77         strcpy(expReg, "v0");
78         return expReg;
79     }
80 }

```

关于形参

由于能力有限，我们仅考虑了单个参数函数和无参数函数两种情况

- 当一个单参数函数调用一个无参数函数，此时栈中只需要存储返回地址寄存器 `$ra`
- 当一个单参数函数调用另一个单参数函数，由于传参寄存器均使用 `$a0`，所以除了返回地址 `$ra` 之外，我们还需要将其压入栈
- 当一个无参数函数调用一个单参数函数，那么这个参数可以由 `move $a0, $<register name>` 挪至函数参数寄存器中，然后再将返回地址寄存器 `$ra` 入栈即可

另外值得一提的是，因为递归函数在调用自身时可能不需要赋值，而是直接利用其形参（如 `int fact(int n): return n*fact(n-1)`），所以我们在处理递归函数调用自身的目标代码时，会特判删去 `move $a0, $<register name>` 这一步传参，直接对 `$a0` 进行操作。

```
1  if(strcmp(head->child->id, currentFunc->child->id)) { // 如果当前翻译的函数不是
    递归的，则需单独传参
2      fprintf(f, "\tmove $a0, $t%d\n", findMark(head->child->bro->bro->child-
    >child->id));
3  }
4
5  if(!strcmp(currentFunc->child->bro->bro->name, "VarList")) { // 如果当前函数是
    有参函数
6      fprintf(f, "\taddi $sp, $sp, -8\n");
7      fprintf(f, "\tsw $a0, 0($sp)\n");
8      fprintf(f, "\tsw $ra, 4($sp)\n");
9      genArgs(head->child->bro->bro, f);
10     fprintf(f, "\tjal %s\n", head->child->id);
11     fprintf(f, "\tlw $a0, 0($sp)\n");
12     fprintf(f, "\tlw $ra, 4($sp)\n");
13     fprintf(f, "\taddi $sp, $sp, 8\n");
14 }
15 else { // 如果当前函数是无参函数
16     fprintf(f, "\taddi $sp, $sp, -4\n");
17     fprintf(f, "\tsw $ra, 0($sp)\n");
18     genArgs(head->child->bro->bro, f);
19     fprintf(f, "\tjal %s\n", head->child->id);
20     fprintf(f, "\tlw $ra, 0($sp)\n");
21     fprintf(f, "\taddi $sp, $sp, 4\n");
22 }
```

实验结果

样例1

```
1 | ./compiler test1.cmm out.s
```

```
1  .data
2  _prompt: .asciiz "Enter an integer:"
3  _ret: .asciiz "\n"
4  .globl main
5  .text
6  read:
7      li $v0, 4
```

```

8      la $a0, _prompt
9      syscall
10     li $v0, 5
11     syscall
12     jr $ra
13 write:
14     li $v0, 1
15     syscall
16     li $v0, 4
17     la $a0, _ret
18     syscall
19     move $v0, $0
20     jr $ra
21 main:
22     li $t1, 0
23     li $t2, 1
24     li $t3, 0
25     addi $sp, $sp, -4
26     sw $ra, 0($sp)
27     jal read
28     lw $ra, 0($sp)
29     addi $sp, $sp, 4
30     move $t4, $v0
31 label1:
32     blt $t3, $t4, label2
33     j label3
34 label2:
35     add $t5, $t1, $t2
36     move $a0, $t2
37     addi $sp, $sp, -4
38     sw $ra, 0($sp)
39     jal write
40     lw $ra, 0($sp)
41     addi $sp, $sp, 4
42     move $t1, $t2
43     move $t2, $t5
44     add $t3, $t3, 1
45     j label1
46 label3:
47     move $v0, $0
48     jr $ra
49

```

样例2

```
1 | ./compiler test2.cmm out.s
```

```

1 | .data
2 | _prompt: .asciiz "Enter an integer:"
3 | _ret: .asciiz "\n"
4 | .globl main
5 | .text
6 | read:
7 |     li $v0, 4

```

```

8      la $a0, _prompt
9      syscall
10     li $v0, 5
11     syscall
12     jr $ra
13 write:
14     li $v0, 1
15     syscall
16     li $v0, 4
17     la $a0, _ret
18     syscall
19     move $v0, $0
20     jr $ra
21 main:
22     addi $sp, $sp, -4
23     sw $ra, 0($sp)
24     jal read
25     lw $ra, 0($sp)
26     addi $sp, $sp, 4
27     move $t1, $v0
28     li $t3, 1
29     bgt $t1, $t3, label4
30     j label5
31 label4:
32     move $a0, $t1
33     addi $sp, $sp, -4
34     sw $ra, 0($sp)
35     jal fact
36     lw $ra, 0($sp)
37     addi $sp, $sp, 4
38     move $t2, $v0
39     j label6
40 label5:
41     li $t2, 1
42 label6:
43     move $a0, $t2
44     addi $sp, $sp, -4
45     sw $ra, 0($sp)
46     jal write
47     lw $ra, 0($sp)
48     addi $sp, $sp, 4
49     move $v0, $0
50     jr $ra
51
52 fact:
53     li $t3, 1
54     beq $a0, $t3, label1
55     j label2
56 label1:
57     move $v0, $a0
58     jr $ra
59     j label3
60 label2:
61     addi $sp, $sp, -8
62     sw $a0, 0($sp)

```



```
63     sw $ra, 4($sp)
64     sub $a0, $a0, 1
65     jal fact
66     lw $a0, 0($sp)
67     lw $ra, 4($sp)
68     addi $sp, $sp, 8
69     mul $v0, $v0, $a0
70     jr $ra
71 label3:
72
```

这里多余的label3仍然合法，不影响运行，故没有进行优化

实验反思

1. 树形IR相对较好的一点在于通过DFS识别模式以生成目标代码，相比线性IR朴素的逐句翻译可以节省操作指令数，更加高效；但相应地，在做寄存器分配优化时就更加不直观，不便于实现优化算法，且因为我们没有设计可以完整存储中间代码的数据结构、而是直接依托于语法树，导致中间代码中ARG、PARAM、CALL等关键字没有派上用场
2. 学期共五次环环相扣的实验，让我深刻理解了管理项目的诸多不易/(T o T)/~~但也着实提高了我的代码实现能力（从画饼到做饼，好不好吃另说）。今后继续加油吧，谢谢师兄和老师的指导！