



南开大学  
Nankai University

南 开 大 学

计 算 机 学 院

编译原理大作业

---

## 从无到有实现编译器

---

年级：2020 级

专业：计算机科学与技术

指导教师：王刚 李忠伟

姓名：马琦

学号：2011600

2023 年 1 月 15 日

## 摘要

本学期，基于编译原理课程所学，我和沈冠翔同学完整实现了从源代码到目标代码的生成。我在本次报告所描述的工作主要如下：首先，借助 lex 词法分析器辅助完成了词法分析的过程，之后我们从产生式入手，设计了完整的产生式，利用 yacc 实现了对产生式的识别。之后为了实现更丰富的功能，我们建立了更多的语法树结点，层次更加丰富，实现了数组的操作。另外我在语法树的建立过程中即进行了类型的检查，实现了对非法数组的检查。之后根据建立的语法树，我们进一步的插入指令，生成对应的中间代码。最后在分配寄存器的时候，我们采用了快速且性能较好的 LinearScan 寄存器分配算法。

**关键字：**词法分析 寄存器分配 语法分析 中间代码生成 类型检查

## 目录

一、 引言	1
二、 编译器整体架构	1
三、 词法分析	2
(一) 正则表达式	2
1. 正则表达式设计	2
2. 举例说明	3
(二) 作用域的确定	3
(三) 函数符号表	3
四、 语法分析	4
(一) 产生式的设计	4
(二) 语法树的建立	4
1. 从整体到局部	4
2. 声明语句	4
3. 赋值式子或者调用	9
4. 表达式	9
5. 函数的声明和调用调用	10
6. While 等语句的处理	10
(三) 数组的类型检查的时机	11
五、 典型的中间代码生成过程	11
(一) Id	12
(二) WhileStmt 和 IfStmt	14
(三) DeclStmt	15
(四) 隐式转换	17
(五) 函数	18
六、 仓库地址	20
七、 寄存器分配	20

**八、 总结**

**20**

## 一、引言

现代场景下的程序开发往往面临着巨量的代码和极为复杂的逻辑，如果单纯使用底层的汇编代码，那么我们开发程序的人工成本将高到我们无法接受，这就是编译器产生的原因，当然前提还是得有一些高级的语言，比如 C++、Java 等，这些语言可以让我们将注意力更多集中在程序本身，而非机器。编译器是为了让机器能够理解我们的高级语言，同时能够完成一些额外的工作，比如能够检查我们的代码是否存在语法错误，甚至可以去优化代码。

本学期基于王刚老师、李忠伟老师讲授的编译原理课程，我们开始从“零”搭建一个简单的基于 SysY 语言的编译器，主要借助了 Yacc、Lex、以及老师学长搭建的架构，能够成功将源代码 (.sy) 翻译生成目标代码 (.s)，在过程中包括对代码的类型检查，同时对中间代码进行了优化，对寄存器的分配进行了优化，使得目标代码的运行效率更高。在这过程中收获很多，具体的细节的一部分将在后文展开。

## 二、编译器整体架构

对于一个编译器，简单来说可以划分为前端、中间代码、后端，其中，前端的目的是将源代码翻译成中间代码，而后端则是将中间代码翻译成机器码。详细来说，当我们传入源代码时，我们首先逐个字得进行识别，生成对应的代词流，将单词流传输给 yacc，进行语法分析，在此过程中我们建立起语法树，同时进行对应的类型检查。之后如果一切正常，我们将利用刚刚生成的语法树进一步地生成中间代码，之后再利用中间代码进一步地生成我们的目标代码。其架构图如下所示。

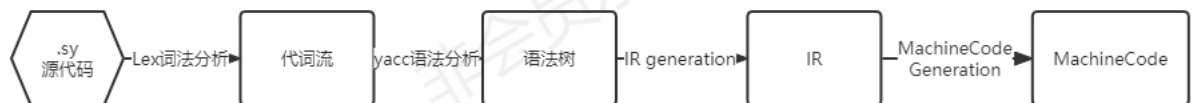


图 1: 朴素的编译器架构

其中，在中间代码到机器码这一个阶段，我们可以加入一些优化来提升代码的性能，首先我们将 IR 转化为 SSA IR (Static Single-Assignment)，之后利用生成的这种 IR 完成 Mem2Reg 和控制流优化算法，之后我们将优化后的 SSA IR 回复为普通的 IR，然后生成机器代码，同时使用 LinearScan 寄存器分配算法对机器代码进行进一步的优化，最后生成我们对应的.s 文件目标代码。

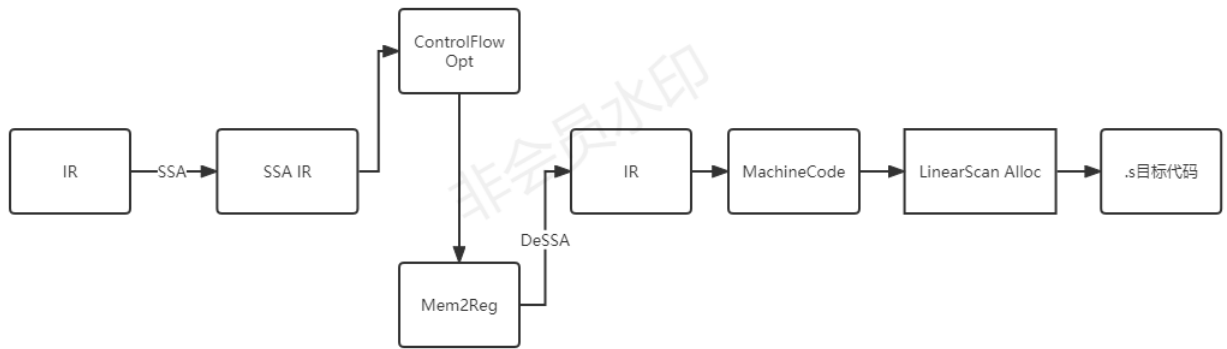


图 2: 带优化的编译器架构

### 三、 词法分析

这一块是编译器中最基础的部分，主要就是对源代码进行切分，将字符按照其类型进行整合，形成合适的单词流供后续分析。这一部分利用到的主要是课上所学的正则表达式、模式匹配，理论工作流程大致如下：

1. Lex 将源文件作为输入，读取源文件的每一个字符，并将它们存储在一个缓冲区中；
2. 然后，Lex 将缓冲区中的字符串与定义的模式（正则表达式）进行比较，如果有匹配的模式，则执行相应的动作；
3. 如果没有匹配的模式，则继续读取源文件，并将读取的字符添加到缓冲区中，重复上述步骤，直到找到匹配的模式。
4. 当找到匹配的模式时，Lex 会执行相应的动作，然后清空缓冲区，重复上述步骤，直到没有更多的字符可以读取。那么我们要做的首先就是设定合适的正则表达式识别不同类型的输入，之后我们还可以利用栈来完成括号匹配并且可以打印出作用域（当然这一部分的内容在后续的 SymbolTable 中有更为完整和强大的功能，当时我们并不知情；，并且这里确实可以打印出作用域）。

#### （一） 正则表达式

##### 1. 正则表达式设计

设计正则表达式需要我们熟悉每一种元素其字符表示类型的特点，例如说对于 ID，第一个字符不可以是数字，对于浮点数来说 e 和 E 只可以出现在最后且并不一定出现，十六进制需要 0x 或者 0X 开头来声明，借助这些特性，我们设计的正则表达式如下：

##### 正则表达式设计

```

1  DECIMAL [1-9][0-9]*|0|-2147483648
2  HEX 0(x|X)[1-9a-fA-F][0-9a-fA-F]*|0
3  OCT 0[1-7][0-7]*
4  DECIMAL_FLOAT ((([0-9]*[.][0-9]*([eE][+-]?[0-9]+)?)|([0-9]+[eE][+-]?[0-9]+)) [fFlLl]? )?
  
```

```

5  HEXADECIMAL_FLOAT (0[xX](((0-9A-Fa-f)*[.][0-9A-Fa-f]*([pP][+-]?[0-9]+)?
   |([0-9A-Fa-f]+[pP][+-]?[0-9]+)))[fLlL]?)
6  ID  [[[:alpha:]]_][[:alpha:]][[:digit:]]_*
7  LINECOMMENT (\\/[^\n]*)
8
9  commentblockbegin  "/*"
10 commentelement  .|\n
11 commentblockend  "*/"

```

## 2. 举例说明

这其中就用到了类似状态机的原理，以一个简单的例子说明，比如要识别 22.33 这个十进制浮点数，首先我们得到 2，此时并不能区分他是浮点数还是十进制整数，但是可以判断他不可能是 ID、十进制和八进制及其他符号，之后识别第二个 2，此时还不能区分，当识别到小数点是，我们可以认为就是浮点数 (如果输入的源文件没有错误)，此时我们继续匹配，当最终匹配到最后的数字时，我们的识别结束，认定该数为十进制浮点数。

## (二) 作用域的确定

这一部分和括号匹配有异曲同工之妙，其实就是如果左大括号入栈之后，右大括号入栈之前 (这里指的是一组匹配的大括号)，那么这之间的所有变量都属于处在作用域。我们只需要用一个全局栈来玩成即可，当然，我们在这里利用 set 加速了查找可用变量的过程 (原来可能需要遍历或者出栈)。由于这一部分在之后的 SymbolTable 中已经有所重复了，这里不做赘述。

作用域的一个例子

```

1  {ID} {
2  if(dump_tokens){
3      if (!cur_str[part].count(std::string(yytext))) cur_str[part].insert(
         std::string(yytext));
4      is_stack[std::string(yytext)].push_back(part);
5      DEBUG_FOR_LEX("ID", std::string(yytext));
6      offset += strlen(yytext);
7  }
8  yylval.strtype = (char*)malloc(strlen(yytext) + 1);
9  strcpy(yylval.strtype, yytext);
10 return ID;
11 }

```

## (三) 函数符号表

在这里需要我们提前将系统库函数加入到符号表里 (函数也被视作 ID 的一种)，因为这些函数在后续不会进行声明定义，因此在类型检查中会被判定为 undefined。这里的操作依然很简单，并且由于我们结合了 multimap，可以在原来的基础上支持一些特性，比如我们可以满足函数和变量同名 (这个功能在 C 语言中有，详细原理在沈冠翔同学的报告中)，下面我们给出一段例子

函数加入符号表

```

1 "getint" {
2     if(dump_tokens){
3         DEBUG_FOR_LEX("ID", "getint");
4         offset += strlen(yytext);
5     }
6     char *lexeme = new char[strlen(yytext) + 1];
7     strcpy(lexeme, yytext);
8     yylval.strtype = lexeme;
9     Type* funcType = new FunctionType(TypeSystem::intType, {}); // 返回类型 int
10    , 无参数
11    SymbolEntry* se = new IdentifierSymbolEntry(funcType, yytext, 0); // 作用域
12    GLOBAL(0)
13    globals->install(yytext, se);
14    return ID;
15 }

```

## 四、 语法分析

### (一) 产生式的设计

这一部分包含了我们编译过程中核心的过程，由之前的词法分析，我们可以得到具有语义的单词流，那么进一步的，我们要根据 SysY 语言的特性，建立对应的语法树，这里我们通过产生式，同样借用模式匹配的思路实现语法的分析，其理论原理为：每次从输入的字符串的最左端开始读取一个字符，然后根据语法规则逐步构建一棵语法树，直到整个字符串都被读取完成，语法树也被构建完成。这里的语法规则其实就是我们之前设计的产生式，在之前的报告中已经详细地给出，这里不做过多赘述。

### (二) 语法树的建立

语法树是我们这个部分需要着重考虑的，后续生成中间代码我们也是按照语法树去生成的。我们在原有的结点里加入一些新的结点，下面我们将按照语法树的生成顺序来介绍我们语法树的建立过程 (以目的为导向)。

#### 1. 从整体到局部

这一部分本可以细分为很多步，但是为了方便理解避免拖沓我们选择集中说明，首先我们把整个源代码当做一个整体，然后不同的代码有不同的块 (也就是有很多作用域)，这其中包含很多不同的语句，包括空语句、赋值语句、条件判断语句、特殊语句 (break、continue、return 等)、循环语句、声明语句、函数定义语句等，我们的语法树也将主要就这些语句展开。

#### 2. 声明语句

我们的程序一般都会操作很多数据，这些数据一般都是声明而来的。所以我们最先开始介绍声明语句 (DeclStmt)，这个语句主要是 "TYPE VARDECLIST;" 或者 "TYPE CONSTVARDECLIST;" 这样的形式，其中 VARDECLIST 有单个定义语句构成，下面我们以非常量的声明语句为例进行说明。

在说明之前，我有必要介绍一下我们新生成的 InitNodeNode 结点，这是一个新的结点，继承自 StmtNode，主要用来存储声明式中的右值，那么相比于原来的 ExprNode，我们加入的这个结点可以满足数组的定义需求，其结构如下：

InitNode 结构

```

1  class InitNode : public StmtNode
2  {
3  private:
4      bool isconst;
5      ExprNode *leaf;
6      std::vector<InitNode *> leaves;
7      int cur_size = 0;
8
9  public:
10     InitNode(bool isconst = false) : isconst(isconst), leaf(nullptr) {};
11     void addleaf(InitNode *next) { leaves.push_back(next); };
12     void setleaf(ExprNode *leaf1) { leaf = leaf1; leaves.clear(); };
13     bool isLeaf() { return leaves.empty(); };
14     void fill(int level, std::vector<int> d, Type* type);
15     int getSize(int d_nxt);
16     bool isFull();
17     bool isConst() const { return isconst; }
18     void output(int level);
19     void genCode(int level);
20     std::vector<InitNode *> getleaves() { return leaves; };
21     ExprNode *getself() { return leaf; };
22     ~InitNode()
23     {
24         if (leaf != nullptr)
25             delete leaf;
26         for (auto _leaf : leaves)
27             delete _leaf;
28     }
29 };

```

这样的一个结点方便我们进一步的递归迭代求解，其逻辑主要如下，当我们是单值或者数组中的单个元素时，leaf 值存储其表达式，而 leaves 为空，如果我们是一个数组时（数组的子数组也算），leaf 为空，leaves 存储其子节点。通过这样的一个简单的逻辑，我们就可以以一种统一的形式表示所有的 Initval。当然对于数组，我们还需要解决的是数组的填充问题，比如下面这一个相对复杂的例子

数组填充示例

```

1  int a[3][2][3] = {{{1}, {2}}, {{3, 4}, 5}};

```

那么对于这样的数组，我们首先需要自顶向下开始遍历，对于外层的大括号，或者就是其第一维度来说，应该有 3 个子节点而此时有两个，所以需要在末尾填充一个结点，如下

数组填充示例



```
1 int a[3][2][3] = {{{1}, {2}}, {{3, 4}, 5}, {}};
```

这其中最简单的是最后一个括号，我们只需要继续按照维度向下遍历，不断加入子节点，最终，他形成

#### 数组填充示例

```
1 {{0, 0, 0}, {0, 0, 0}}
```

而对于第二种来说,我们直接判断叶节点的数量肯定是不行的,我们首先需要倒着遍历一遍 leaves, 数一下尾部为单个元素的叶节点, 之后将其补成可以被当前维度容量整除的数量, 之后再按照下一维度的要求挨个递归插入, 其过程如下

#### 数组填充示例

```
1 {{3, 4}, 5, 0}
2 {{3, 4}, 5, 0, {}}
3 {{3, 4}, 5, 0, {0, 0}}
```

第一种也不是很麻烦, 我们直接递归, 当递归到数组的最后一个维度时, 选择填充 0 元素, 并且退出递归, 其过程如下

#### 数组填充示例

```
1 {{1}, {2}, {}}
2 {{1, 0}, {2, 0}, {0, 0}}
```

这一部分对应的代码如下

#### 数组填充代码

```
1 void InitNode::fill(int level, std::vector<int> d, Type *type) // 传入维度信
   息和类型, 并且记录递归深度
2 {
3     if (level == d.size() || leaf != nullptr)
4     {
5         if (leaf == nullptr) // 递归终点
6         {
7             setleaf(new Constant(new ConstantSymbolEntry(Var2Const(type), 0))
8                 );
9         }
10        return;
11    }
12    int cap = 1, num = 0, cap2 = 1;
13    for (int i = level + 1; i < d.size(); i++) // 计算数组容量
14        cap *= d[i];
15    for (int i = leaves.size() - 1; i >= 0; i--) // 计算尾部单个元素数量
16        if (leaves[i]->isLeaf())
17            num++;
18    else
19        break;
20    while (num % cap) // 填充成为容量的整数倍
```

```

21     {
22         InitNode *new_const_node = new InitNode(true);
23         new_const_node->setleaf(new Constant(new ConstantSymbolEntry(
24             Var2Const(type), 0)));
25         addleaf(new_const_node);
26         num++;
27     }
28     int t = getSize(cap); // 补充数组结点
29     while (t < d[level])
30     {
31         InitNode *new_node = new InitNode(true);
32         addleaf(new_node);
33         t++;
34     }
35     for (auto l : leaves) // 进一步递归
36     {
37         l->fill(level + 1, d, type);
38     }
39 int InitNode::getSize(int d_nxt) // 计算当前结点下子节点的数量(非单个结点)
40 {
41     int num = 0, cur_fit = 0;
42     for (auto l : leaves)
43     {
44         if (l->leaf != nullptr)
45         {
46             num++;
47         }
48         else
49         {
50             cur_fit++;
51         }
52         if (num == d_nxt)
53         {
54             cur_fit++;
55             num = 0;
56         }
57     }
58     return cur_fit + num / d_nxt;
59 }

```

填充的时机在 DeclStmt 声明时即可，这样可以让我们尽早获得规整的数组结构。下面开始介绍声明式子，针对普通的整型或者浮点型变量来说，声明语句一般像下面这样

#### 声明语句示例

```

1 int a = 1;
2 float b = 1.0;

```

我们首先在 DeclStmt 的级别 (即可以得到 TYPE 时) 先将类型存储起来，之后识别到 ID 时，

首先进行类型检查，判断之前是否重复定义过，之后我们将类型存入符号表中，并根据这个 SymbolEntry 生成一个新的 Id 结点，Id 结构主要如下

Id 结点的结构

```

1  class Id : public ExprNode
2  {
3  private:
4      IndicesNode *indices; //用来存储Id的索引，只对数组有用
5      bool is_array = false, is_array_ele = false; // 用来区分数组还是数组元素
6      bool isleft; // 判断是否为左值
7
8  public:
9      Id(SymbolEntry *se, bool be_array = false, bool isleft = false) :
10         ExprNode(se, be_array)
11     {
12         indices = nullptr;
13         is_array_ele = se->getType()->isArray() && be_array;
14         is_array = se->getType()->isArray();
15     };
16     void setIndices(IndicesNode *new_indices) { indices = new_indices; };
17     IndicesNode *getIndices() { return indices; };
18     void output(int level);
19     bool isArray() { return is_array; };
20     bool isArrayEle() { return is_array_ele; };
21     ArrayType* getArray_Type() { return (ArrayType*)(getSymPtr()->getType())
22         ; };
23     // void typeCheck();
24     void genCode();
25     ~Id()
26     {
27         if (indices != nullptr)
28             delete indices;
29     };
30 };

```

之后将该 Id 存到符号表中，再将 Initval 所生成的 InitNode 和当前 Id 一起，通过如下语句生成一个新的声明语句结点。

Id 结点的结构

```

1  $$ = new DeclStmt(new Id(se), (dynamic_cast<InitNode*>($3)));

```

对于数组的声明式来说，是一样的道理，主要区别在于数组有维度信息，所以我们又新定义了一个 IndicesNode 结点，结点的结构如下

IndicesNode 结构

```

1  private:
2      std::vector<ExprNode *> exprList; //存储维度信息
3
4  public:

```

```

5     IndicesNode() {};
6     void addNew(ExprNode *new_expr) { exprList.push_back(new_expr); };
7     void addBefore(ExprNode *new_expr) { exprList.insert(exprList.begin(),
8         new_expr); };
9     void output(int level);
10    std::vector<ExprNode *> getExprList() { return exprList; };
11    void genCode();
12    ~IndicesNode()
13    {
14        for (auto expr : exprList)
15            delete expr;
16    }

```

在数组维度的产生式中，我们每添加一个 Expr(或者 ConstExpr) 直接将其加入到数组容器的最后就可以了，最后我们依然是在 DeclStmt 层次，先将 IndicesNode 值赋给 Id，之后重复之前普通值的操作即可。

### 3. 赋值式子或者调用

按照如下这种情况说明

#### 赋值及调用示例

```

1 int a = 1;
2 int b[2][2] = {a};

```

这里我们应该把第二个“a” 识别为 LRval，也就是已经定义过的值，这样就可以和定义时的 ID 区分开了。数组同样可行，添加维度信息即可。

### 4. 表达式

这虽然是我们很早就完成的部分，但是确实后期一直会频繁调用的部分，这里主要包括我们的 ExprNode，这主要从 Exp(或者 ConstExp) 开始算起，之后自顶向下最先到的应该是优先级最低的加减法表达式，即 ()+() 或者 ()，之后每一个都可以继续迭代出更多的加减法表达式，之后进入下一级乘除表达式，主要包含乘法除法和取模运算，也可以继续迭代，当然也可以就此收手，进入 UnaryExpr，也就是元表达式，可以理解为无法继续拆分的运算单位，包括之前的 Id、常数等等。当然还有一种表达式，即 if、while 语句中的判断表达式，这一部分和之前的式子差别不大，主要在于最后需要添加隐式转换结点，隐式转换结点如下：

#### 隐式转换结点

```

1 class ImplicitCast : public ExprNode
2 {
3     private:
4         ExprNode *expr;
5
6     public:
7         ImplicitCast(SymbolEntry *se, ExprNode *expr) : ExprNode(se), expr(expr)
8         {};
9         void output(int level);
10        // void typeCheck();

```

```

10 void genCode();
11 ~ImplicitCast() { delete expr; };
12 };

```

之后我们只需要在生成中间代码的时候将类型转换即可。

## 5. 函数的声明和调用

我们在这里实现了一个新的结点 FuncDefNode, 其结构如下

函数声明结点

```

1 class FuncDefNode : public StmtNode
2 {
3 private:
4     SymbolEntry *se; // 存符号表
5     FuncDefParamsNode *params; // 存变量, 这里节省地方就不写了
6     StmtNode *stmt; // 存表达式
7
8 public:
9     FuncDefNode(SymbolEntry *se, FuncDefParamsNode *params, StmtNode *stmt,
10                 bool needRet) : se(se), params(params), stmt(stmt) {};
11     void output(int level);
12     // void typeCheck();
13     void genCode();
14     ~FuncDefNode()
15     {
16         if (params != nullptr)
17             delete params;
18         delete stmt;
19     };
20 };

```

我们在定义时获得的是 Type FuncName(FuncParamsList) Stmts, 此时我们先声明一个新的 Id 结点, 并且将函数的参数列表及其声明式 (BlockStmt) 和 Id 一起生成一个新的 FuncDefNode 结点, 如下

函数定义

```

1 $$ = new FuncDefNode(se, dynamic_cast<FuncDefParamsNode*>($5), $8, needRet);

```

值得注意的是, 其中 needRet 存储的是函数是否需要函数返回值, 特殊标记即可。在函数调用时, 和普通 Id 的区别在于添加了参数, 其余并无区别。

函数调用的参数中如果出现数组, 那么我们只需要额外地标记一下即可, 比如把第一个空维度标记为-1, 具体的区别我们会在中间代码的部分进行描述。

## 6. While 等语句的处理

这一部分在语法分析并不难, 我们只需要记录当前 While(if) 条件判断的式子以及语句块内的表达式, 在之后的中间代码生成部分才是重点。

### (三) 数组的类型检查的时机

数组检查有以下几个时机:

- 数组填充: 当我们补齐末尾后 (详见之前数组填充), 需要继续补充下一维度的结点, 如果此时已经超过当前维度的容量那么直接报错。

错误数组

```
1      int a[4][2] = {{1}, {2}, 3, 4, 5, 6, 7, 8, 9};
2      // fill -> {{1}, {2}, 3, 4, 5, 6, 7, 8, 9, 0}
3      // size = 2 + (7 + 1) / 2 = 6 > 4
4      // quit
```

- 访问数组结点: 如果当前的访问超过了维度限制, 直接报错

错误数组

```
1      int a[4][2] = {{1}, {2}, 3, 4, 5, 6, 7, 8, 9};
2      a[5][1];
3      // 5 > 4
4      // quit
```

- 赋初值时: 声明新的 DeclStmt 是, 应该判断第一层 (不递归) InitNode 里面的 leaf 是否不为空, 不为空则报错。

错误数组

```
1      int a[4][2] = 1;
2      // InitNode(leaf = ExprNode(1), leaves.size() == 0)
3      // quit
```

- 数组填充计数: 在统计数组当前维度的叶节点时, 我们存了两个, 一个是 cur\_fit, 一个是 num, 其中 num 统计的是单个元素的值, 当 cur\_fit 增加时, 如果 num 不为 0, 则返回报错

错误数组

```
1      int a[4][2] = {{1, 2}, 3, {4}};
2      // 当到{4}时, cur_fit应该加一, 但是此时num为1, 所以说明之前的单个
      // 元素并没有构成当前维度的一块
3      // quit
```

## 五、典型的中间代码生成过程

这一块主要是结合我们之前语法树的部分, 从 root 开始, 自顶向下逐层生成中间代码, 至于语法树的生成逻辑我们之前已经谈过了, 这里不再赘述。下面只介绍一些比较特别的中间代码生成。

## (一) Id

Id 的中间代码需要区分其是数组元素还是普通元素，主要区别在于取指操作，那么仿照我们第一次实验模拟 LLVM 语言的经历，此处我们应该使用类似于 `getelementptr` 指令进行偏移，这里唯一需要注意的就是，我们每次寻址时，指针的类型会不断发生变化，即从 `[]...(n 个)` 到 `[]...(n-1 个)`，所以我们需要不断地修改操作数的类型，具体的代码如下：

## IdgenCode

```

1 void Id::genCode()
2 {
3     if (getType()->isConst() && !is_Array()) // 常量折叠
4         return;
5     BasicBlock *bb = builder->getInsertBB();
6     Operand *addr = dynamic_cast<IdentifierSymbolEntry *>(symbolEntry)->
7         getAddr();
8     if (!is_Array() && !is_Array_Ele()) // 普通的genCode
9     {
10         // delete dst;
11         dst = new Operand(new TemporarySymbolEntry(symbolEntry->getType(),
12             SymbolTable::getLabel()));
13         new LoadInstruction(dst, addr, bb);
14     }
15     else
16     {
17         SymbolEntry *temp = new TemporarySymbolEntry(new PointerType(((
18             ArrayType *)getType())->getElemType()), SymbolTable::getLabel());
19         // 获得元素类型
20         if (indices != nullptr)
21         {
22             Operand *tempSrc = addr;
23             ArrayType* curr_type; // 生成新的类型
24             if (curr_type->isIntArray()) {
25                 if (curr_type->isConst())
26                     curr_type = new ConstIntArrayType();
27                 else curr_type = new IntArrayType();
28             }
29             else {
30                 if (curr_type->isConst())
31                     curr_type = new ConstFloatArrayType();
32                 else curr_type = new FloatArrayType();
33             }
34             curr_type->SetDim(get_Array_Type()->fetch());
35             std::vector<int> curr_dim = curr_type->fetch();
36             Operand *tempDst = new Operand(new TemporarySymbolEntry(new
37                 PointerType(curr_type), SymbolTable::getLabel()));
38             // fprintf(stderr, "id type size is %d, idx size is %d\n",
39                 curr_dim.size(), indices->getExprList().size());
40             bool is_first = true;

```

```

35     for (auto idx : indices->getExprList())
36     {
37         if (!is_first) {
38             fprintf(stderr, "last_op_type is %s\n", last_op->getType
39                 ()->toStr().c_str());
40         }
41         idx->genCode();
42         auto gep = new GepInstruction(tempDst, tempSrc, idx->
43             getOperand(), bb, flag);
44         // 生成跳转指令
45         last_op = tempSrc;
46         if (!curr_dim.empty()) {
47             // 更新类型的维度信息
48             curr_dim.erase(curr_dim.begin());
49             curr_type->SetDim(curr_dim);
50         }
51         else break;
52         tempSrc = new Operand(new TemporarySymbolEntry(new
53             PointerType(curr_type), ((TemporarySymbolEntry*)tempDst->
54                 getEntry())->getLabel()));
55         tempDst = new Operand(new TemporarySymbolEntry(new
56             PointerType(curr_type), SymbolTable::getLabel()));
57         is_first = false;
58     }
59     Operand* new_dst = new Operand(new TemporarySymbolEntry(new
60         PointerType(curr_type->getElemType()), ((TemporarySymbolEntry
61             *)tempDst->getEntry())->getLabel()));
62     if (!pointer)
63     {
64         // 数组元素要加入load指令(如果需要用到值的话)
65         Operand* dst1;
66         if (getType()->isInt())
67             dst1 = new Operand(new TemporarySymbolEntry(TypeSystem::
68                 intType, SymbolTable::getLabel()));
69         else dst1 = new Operand(new TemporarySymbolEntry(TypeSystem::
70             floatType, SymbolTable::getLabel()));
71         new LoadInstruction(dst1, new_dst, bb);
72         dst = dst1;
73     }
74 }
75
76 ....
77
78 }
79
80 }

```



## (二) WhileStmt 和 IfStmt

这一部分比较重要的是回填 (BackPatch), 我们需要利用栈来存储上一个语句块的返回地址, 以此来记录保存返回地址等信息, 回填本质上是一种延时更新, 当我们不知道当前的 true 和 false 分别返回哪里时, 我们可以先搁置, 等确认后再回来更新。另外这里的条件还涉及逻辑短路的问题, 这个也不是很困难, 我们只需要按照顺序, 分别设置表达式的 true 和 falselist 即可, 相当于有这么一种结构。

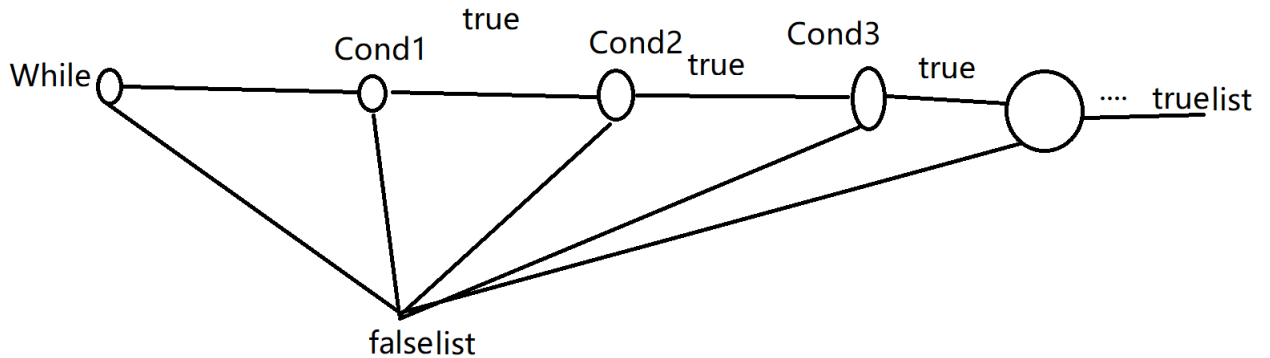


图 3: 短路机制

这里的代码实现起来相当简单, 我们以 IfElseStmt 举例

### IfElseStmtgenCode

```

1 void IfElseStmt::genCode()
2 {
3     Function *func = builder->getInsertBB()->getParent();
4     BasicBlock *then_bb, *else_bb, *end_bb = new BasicBlock(func);
5
6     height = 1;
7     cond->genCode();
8     height = 0;
9
10    if (!cond->trueList().empty()) // 常量折叠, 排除 if(0)
11    {
12        then_bb = new BasicBlock(func);
13        backPatch(cond->trueList(), then_bb);
14        builder->setInsertBB(then_bb);
15        thenStmt->genCode();
16        then_bb = builder->getInsertBB();
17        new UncondBrInstruction(end_bb, then_bb);
18    }
19
20    if (!cond->>falseList().empty()) // 常量折叠, 排除 if(1)
21    {
22        else_bb = new BasicBlock(func);
23        backPatch(cond->>falseList(), else_bb);
24        builder->setInsertBB(else_bb);

```

```

25     elseStmt->genCode();
26     else_bb = builder->getInsertBB();
27     new UncondBrInstruction(end_bb, else_bb);
28 }
29
30 builder->setInsertBB(end_bb);
31 }

```

另外值得注意的是，对于 While 的返回地址，我们需要考虑 While() While() .... 这种嵌套情况，所以应该用栈来模拟作用域，然后只要将 falselist 设置为栈顶元素即可。

### (三) DeclStmt

这一部分是修改最多的一部分，从初始化开始就充满了戏剧性。

#### DeclStmtInit

```

1 DeclStmt(Id *id, InitNode *expr = nullptr, bool isConst = false, bool isArray
  = false) : id(id), expr(expr), BeConst(isConst), BeArray(isArray)
2 {
3     next = nullptr;
4     if (expr != nullptr) {
5         fprintf(stderr, "-----\n");
6         if (id->getType()->isArray()) {
7             std::vector<int> origin_dim = ((ArrayType*)(id->getType()))->
              fetch();
8             expr->fill(0, origin_dim, ((ArrayType*)(id->getType()))->
              getElemType());
9         }
10    }
11 };

```

这里我们需要填充 initNode, 这是为了方便后续初始化。在生成中间代码的部分里，我们主要用来区分全局变量和局部变量，全局变量变量比较简单，我们主要说一下局部变量。对于普通的整型或者浮点型，我们找个位置 Store 就行了，而对于数组来说，我们不仅要找个位置，还要把他后面的所有值都 Store，所以这里，我们依然采用 Gep 指令，先迭代几轮，然后到那个位置 Store 就可以了，值得注意的是，直接拿 InitNode 去初始化非常得不方便，因为你不知道他是哪个位置 (比如 1, 0, 2, 0, 3, 4, 5, 6 里面的 4，直观你不知道他在哪，递归知道了第二层，然后是第三个结点的唯一元素，难不成是 [3][0])，所以我们这里来个 DFS 把所有的叶节点全部先拿出来存起来，具体的代码如下：

#### Preprocess

```

1 std::vector<ExprNode *> vec_val;
2 static void get_vec_val(InitNode *cur_node)
3 {
4     if (cur_node->isLeaf() || cur_node->getself() != nullptr)
5     {
6         vec_val.push_back(cur_node->getself());
7     }

```

```

8     else
9     {
10         for (auto l : cur_node->getleaves())
11         {
12             get_vec_val(l);
13         }
14     }
15 }

```

这样我们就将一个多维数组展开了，那么具体填充时，由计算公式

$$Addr_{arr[k_1]...[k_n]} = k_1 \times \prod_{i=2}^n N_i + \dots + k_i \times \prod_{j=i+1}^n N_j + \dots k_n \times 1$$

其中 N 为 arr 的形状，那么我们每次只需要除一下最大的数就行了，其他数除完了都是 0，这样递归就可以算出最终的地址了，具体的代码如下

#### InitNodegenCode

```

1 void InitNode::genCode(int level)
2 {
3     for (int i = 0; i < vec_val.size(); i++)
4     {
5         int pos = i;
6         std::vector<int> curr_dim(curr_dim);
7         Operand *final_offset = arrayAddr;
8         for (int j = 0; j < d.size(); j++)
9         {
10             cur_type->SetDim(curr_dim); // 更改形状，也可以理解成一个不断减去
                第一个维度的数组类型
11             curr_dim.erase(curr_dim.begin());
12             Operand *offset_operand = new Operand(new ConstantSymbolEntry(
                TypeSystem::constIntType, pos / d[j]));
13             Operand *addr = final_offset;
14             final_offset = new Operand(new TemporarySymbolEntry(new
                PointerType(cur_type), SymbolTable::getLabel()));
15             fprintf(stderr, "currdim is %s\n", addr->getType()->toStr().c_str
                ());
16             new GepInstruction(final_offset, addr, offset_operand, builder->
                getInsertBB()); //跳一下，跳到没超过目标的最大位置
17             pos %= d[j]; // 每一次求一下余数就是剩下的
18         }
19         vec_val[i]->genCode();
20         Operand *src = vec_val[i]->getOperand();
21         final_offset = new Operand(new TemporarySymbolEntry(
22             new PointerType(((ArrayType *)cur_type)->getElemType()),
23             dynamic_cast<TemporarySymbolEntry *>(final_offset->getEntry()->
                getLabel()))); //要修改一下，此时的指针类型不应该再是数组，我们
                把操作数取出来，改一下类型就好了
24         new StoreInstruction(final_offset, src, builder->getInsertBB()); //找
                到最后的位置直接store

```

```

25     curr_dim.clear();
26     // assert(cur_dim.empty());
27 }
28 }

```

#### (四) 隐式转换

这个之前已经说过，我们在这里直接贴代码，主要就是把操作数的类型修改一下即可

##### 隐式转换

```

1 void ImplicitCast::genCode()
2 {
3     BasicBlock *bb = builder->getInsertBB();
4     Function *func = (bb == nullptr) ? nullptr : bb->getParent();
5     if (getType()->isConst()) // 常量折叠
6     {
7         if (getType() == TypeSystem::constBoolType)
8             if (height > 0)
9             {
10                 if (getValue())
11                 {
12                     BasicBlock *trueBB = nullptr;
13                     true_list.push_back(new UncondBrInstruction(trueBB, bb));
14                 }
15                 else
16                 {
17                     BasicBlock *falseBB = nullptr;
18                     false_list.push_back(new UncondBrInstruction(falseBB, bb));
19                 }
20             }
21         return;
22     }
23     expr->genCode();
24     // int/float -> bool
25     if (this->getType()->isBool() || this->getType()->isConstBool())
26     {
27         Operand *src1 = expr->getOperand(), *src2 = new Operand(new
                ConstantSymbolEntry(Var2Const(src1->getType()), 0));
28         new CmpInstruction(CmpInstruction::NE, dst, src1, src2, bb);
29         if (height > 0)
30         {
31             BasicBlock *trueBB = nullptr, *falseBB = new BasicBlock(func), *
                endBB = nullptr;
32             true_list.push_back(new CondBrInstruction(trueBB, falseBB, dst,
                bb));
33             false_list.push_back(new UncondBrInstruction(endBB, falseBB));
34         }

```

```

35     }
36     else if (this->getType()->isInt() || this->getType()->isConstInt())
37     {
38         Operand *src = expr->getOperand();
39         // bool -> int
40         if (src->getType()->isBool() || src->getType()->isConstBool())
41             new ZextInstruction(dst, src, bb);
42         // float -> int
43         else
44         {
45             assert(src->getType()->isFloat() || src->getType()->isConstFloat());
46             new IntFloatCastInstruction(IntFloatCastInstruction::F2S, dst,
47                 src, bb);
48         }
49     }
50     else
51     {
52         assert(this->getType()->isFloat() || this->getType()->isConstFloat());
53         ;
54         Operand *src = expr->getOperand();
55         // bool -> float
56         if (src->getType()->isBool() || src->getType()->isConstBool())
57         {
58             Type *type = src->getType()->isConstBool() ? TypeSystem::
59                 constIntType : TypeSystem::intType;
60             Operand *t = new Operand(new TemporarySymbolEntry(type,
61                 SymbolTable::getLabel()));
62             new ZextInstruction(t, src, bb);
63             new IntFloatCastInstruction(IntFloatCastInstruction::S2F, dst, t,
64                 bb);
65         }
66         // int -> float
67         else
68         {
69             assert(src->getType() == TypeSystem::intType || src->getType() ==
70                 TypeSystem::constIntType);
71             new IntFloatCastInstruction(IntFloatCastInstruction::S2F, dst,
72                 src, bb);
73         }
74     }
75 }

```

## (五) 函数

函数声明可以看作包含其他所有的大块，我们只需要把其中的所有分支之间的关系弄清楚就可以了，不多说了，上代码。

## 函数

```

1 void FuncDefNode::genCode()
2 {
3     Unit *unit = builder->getUnit();
4     Function *func = new Function(unit, se);
5     BasicBlock *entry = func->getEntry();
6     // set the insert point to the entry basicblock of this function.
7     builder->setInsertBB(entry);
8     if (params != nullptr)
9         params->genCode();
10    stmt->genCode();
11    for (auto bb = func->begin(); bb != func->end(); bb++)
12    {
13        for (auto i : (*bb)->begin(); i != (*bb)->end(); i = i->getNext())
14            if (i->isRet())
15            {
16                while (i != (*bb)->rbegin())
17                    delete (index->getNext());
18                break;
19            }
20        // 获取一下最后一块指令
21        Instruction *last = (*bb)->rbegin();
22        // 把带条件的修改了
23        if (last->isCond())
24        {
25            BasicBlock *trueBB = dynamic_cast<CondBrInstruction *>(last)->
26                getTrueBranch();
27            BasicBlock *falseBB = dynamic_cast<CondBrInstruction *>(last)->
28                getFalseBranch();
29            (*bb)->addSucc(trueBB);
30            (*bb)->addSucc(falseBB);
31            trueBB->addPred(*bb);
32            falseBB->addPred(*bb);
33        }
34        // 无条件的直接加在后边就行了
35        if (last->isUncond())
36        {
37            BasicBlock *dstBB = dynamic_cast<UncondBrInstruction *>(last)->
38                getBranch();
39            (*bb)->addSucc(dstBB);
40            dstBB->addPred(*bb);
41        }
42    }
43 }

```

其他的函数参数调用就不说了，具体的看代码吧。

## 六、 仓库地址

认准`nku_stack`。

## 七、 寄存器分配

我们在这里实现了一个低配版的最优调度问题，大致的流程如下：在编译器中，首先从源代码生成中间代码，然后按照程序的指令顺序从上到下扫描中间代码，根据每条指令的寄存器使用情况，把寄存器分配给每条指令，最后把中间代码转换为机器码。这个使用情况是这样进行描述的，我们把每一个寄存器活跃的时间段按照起始时间排个序，然后对于使用时间段不冲突的可以考虑用相同的寄存器，如果分配的寄存器超过了我们能提供的，直接放到 Mem 里面，这就是 Spill 操作，这里我们优先考虑 Spill 距离当前最远的，也是一种启发手段 (难道要联动体系结构的 replacement 吗:)?)。

这里的代码太长，不宜体现主旨，详细内容可以进仓库查看。

## 八、 总结

经过这一个学期的努力，和队友终于完成了这么大的一个项目，加起来应该有几千行代码，其中很多的思路也并非简单的模拟，需要我们熟悉原理的同时设计算法去优化，Debug 的过程痛苦但是有趣 ()，过程中收获颇丰。

**感谢老师的悉心教导!**

**感谢学长学姐热心答疑!**

**也感谢我的神级队友一路相随、不离不弃带我飞!!!!!!**

最后感谢编译原理让我认识到原来 C++ 代码也可以这么难写。