

编译原理实验报告

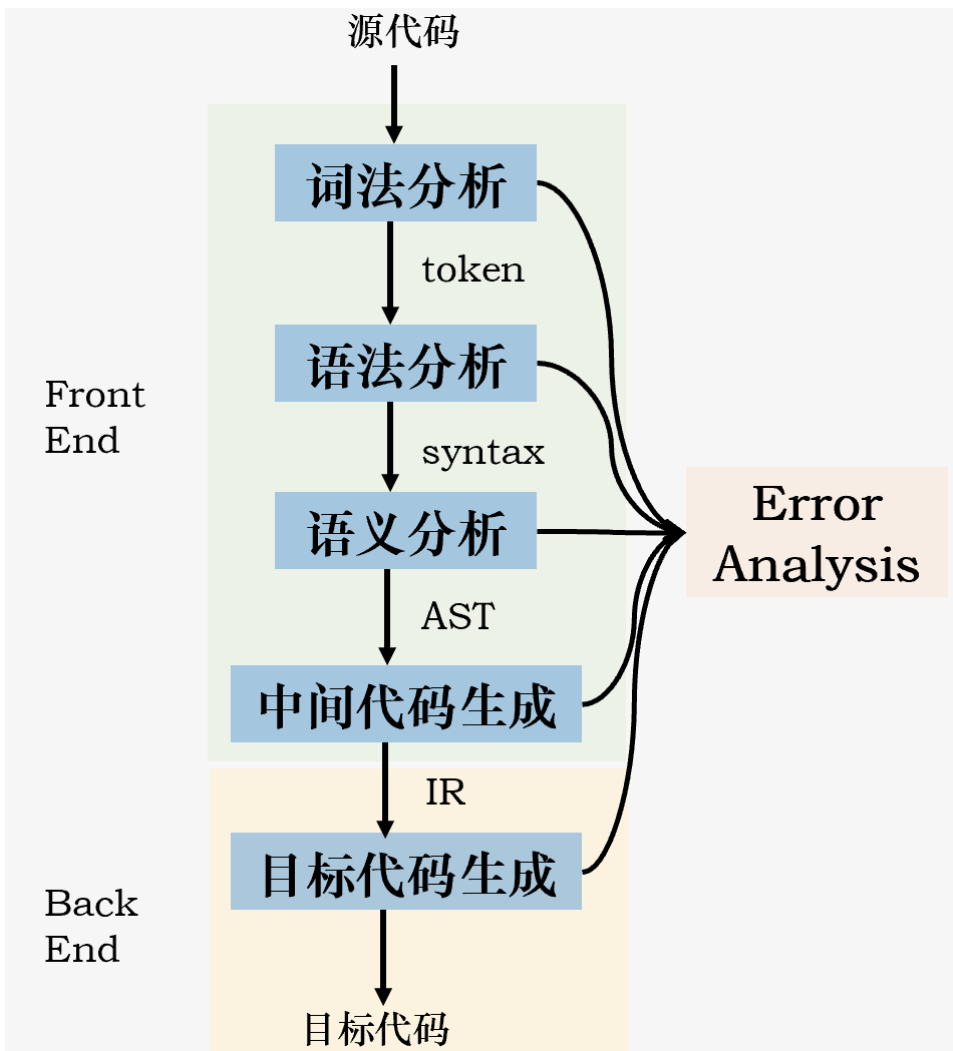
C-MINUS Compiler

3200104880 王晶晶

2023.5.28

引言

C-MINUS Compiler 主要由五个阶段组成：**词法分析**、**语法分析**、**语义分析**、**中间代码生成**、**目标代码生成**。其中中间代码生成之前属于前端，目标代码生成属于后端。词法分析使用Flex将定义的正则表达式转换成对应的状态机，语法分析使用Bison将定义的文法转换成对应的parser，语义分析为构建AST树，调用Illum库来生成中间代码，并借助clang编译生成目标代码。



运行环境

语言

- C++

工具链

- cmake
- flex, bison, llvm-14

实验环境

- windows amd64
- wsl2 + ubuntu 20.04

词法分析

Flex

- Flex是一个生成词法分析器的工具，它可以利用正则表达式来生成匹配相应字符串的C语言代码，其语法格式基本同Lex相同。
- 单词的描述称为模式(Lexical Pattern)，模式一般用正规表达式进行精确描述。Flex通过读取一个有规定格式的文本文件，输出一个C语言源程序。
- Flex的输入文件为lex源文件，它内含正规表达式和对相应模式处理的C语言代码。lex源文件的扩展名习惯上用.l表示。Flex通过对源文件的扫描自动生成相应的词法分析函数int yylex()，并将之输出到名规定为lex.yy.c的文件中。

CMINUS Lexer

- tokens
 - `int, double, char`
声明类型，CMINUS编译器支持三种数据类型，int表示32位整型，double表示浮点型，char表示8位字符型。
 - `*, /, +, -, %`
算数运算，CMINUS编译器支持五种基础运算，包括以上三种数据类型的乘法、除法、加法、减法、取模运算。
 - `>, <, ==`
比较运算，CMINUS编译器支持三种基础比较运算，包括大于、小于、等于。
 - `=`
赋值运算符。
 - `&`
取地址运算符。
 - `(,), [,], {, }, ;`
CMINUS基本分隔符。`()`表示基本expression, `[]`用于数组元素访问, `{ }`表示基本block, 用于分隔一个新的作用域。
`;`为句末提示符。
 - `if, else, while, return`
CMINUS逻辑语句关键词，支持if, else 条件语句, while 循环语句, return 函数返回。

Lexer 实现

```
1  "="                                { std::cout<<"ASSIGN"<<std::endl; return
    TOKEN(ASSIGN); }
2  "=="                               return TOKEN(EQ);
3  "!="                               return TOKEN(NEQ);
4  "<"                                return TOKEN(LT);
5  "<="                              return TOKEN(LE);
6  ">"                                return TOKEN(GT);
7  ">="                              return TOKEN(GE);
8  "("                                { std::cout<<"LPAREN"<<std::endl; return
    TOKEN(LPAREN); }
9  ")"                                { std::cout<<"RPAREN"<<std::endl; return
    TOKEN(RPAREN); }
10 "["                               return TOKEN(LBRACKET);
11 "]"                               return TOKEN(RBRACKET);
12 "{"                                { std::cout<<"LBRACE"<<std::endl; return
    TOKEN(LBRACE); }
13 "}"                                { std::cout<<"RBRACE"<<std::endl; return
    TOKEN(RBRACE); }
14 ","                               return TOKEN(COMMA);
15 ";"                               { std::cout<<"COMMA"<<std::endl; return
    TOKEN(SEMI); }
16 "+"                               return TOKEN(PLUS);
17 "-"                               return TOKEN(MINUS);
18 "*"                               return TOKEN(MUL);
19 "/"                               return TOKEN(DIV);
20 "%"                               return TOKEN(MOD);
21 "&&"                              return TOKEN(AND);
22 "||"                              return TOKEN(OR);
23 "&"                                return TOKEN(ADDR);
24 "if"                              return TOKEN(IF);
25 "else"                            return TOKEN(ELSE);
26 "while"                           return TOKEN(WHILE);
27 "return"                          { std::cout<<"return"<<std::endl; return
    TOKEN(RETURN); }
28 "break"                           return TOKEN(BREAK);
29 [a-zA-Z_][a-zA-Z0-9_]*           {SAVE_TOKEN; std::cout<<yylval.string<<std::endl;
    return IDENTIFIER; }
30 [0-9]+                             SAVE_TOKEN; return INTEGER;
31 [0-9]+\.[0-9]*                     SAVE_TOKEN; return DOUBLE;
32 ["'].*["']                         SAVE_TOKEN; return STRING;
33 "\'\"[^\\"']\"'"                  SAVE_TOKEN; return CHARACTER;
34 [ \t\n]                            { ; }
35 .                                  { printf("Unknown token!\n");
    std::cout<<yylval.string<<std::endl; }
```

语法分析

Bison

- Bison通过提供的产生式多次构造，最终得到一个动作表，然后利用这个动作表去解析句子。
- bison读取用户提供的语法产生式，生成一个C语言格式的动作表，并将其包含进一个名为 `yyparse()` 的C函数，这个函数的作用是利用这个动作表来解析token流，这个token流是由flex生成的词法分析器扫描源程序得到。

Grammar

主要分为以下几个产生式，来生成目标CMINUS语法表达。`program`是起始符(start symbol)，所有的程序都由`program`开始生成。`stmts`是由多个`stmt`组合而成的，`stmt`是语言中的基本语句块，包含变量声明、函数声明、表达式、返回语句、条件语句、循环语句等多种语句。`block`表示由`{}`包裹的基本变量作用域。`expr`表示基本表达式，如算术运算、变量赋值等等。

- `program`

```
1 program : stmts {
2         root = $1;
3         std::printf("begin parsing!\n");
4     };
```

`program`作为程序的根节点，用于生成多个`stmts`。

- `stmts`

```
1 stmts : stmt {
2         $$ = new Block();
3         $$->statements.push_back(<stmt_type>1);
4         std::printf("stmt\n");
5     } |
6     stmts stmt {
7         $1->statements.push_back(<stmt_type>2);
8     };
```

`stmts`初始化一个`Block`节点，用于存储多个`stmt`。

- `stmt`

`stmt`包含变量声明

```
1 stmt : var_decl SEMI {
2         std::printf("vardecl\n");
3     } |
```

函数声明及定义

```
1 stmt : func_decl {
2         std::printf("funcdecl\n");
3     } |
```

表达式

```
1 stmt : expr SEMI {
2         $$ = new ExpressionStatement(*$1);
3     } |
```

`void`型和含参返回语句

```

1 stmt : RETURN SEMI {
2         $$ = new ReturnStatement();
3         std::printf("return\n");
4     }|
5     RETURN expr SEMI {
6         $$ = new ReturnStatement($2);
7         std::printf("return\n");
8     }|

```

break语句

```

1 stmt : BREAK SEMI {
2         $$ = new BreakStatement();
3     }|

```

条件语句

```

1 stmt : IF LPAREN expr RPAREN block ELSE block{
2         $$ = new IfStatement($3, $5, $7);
3     }|
4     IF LPAREN expr RPAREN block{
5         $$ = new IfStatement($3, $5);
6     }|

```

循环语句

```

1 expr : WHILE LPAREN expr RPAREN block {
2         $$ = new WhileStatement($3, $5);
3     };

```

- block

block 包含由 {} 包裹起来的变量作用域。

```

1 block : LBRACE stmts RBACE {
2         $$ = $2;
3         std::printf("block!\n");
4     }|
5     LBRACE RBACE{
6         $$ = new Block();
7     };

```

- expr

expr可以生成多种基本表达式：变量赋值、变量取地址、立即数初始化、算术运算、逻辑运算等

```

1 expr : identifier ASSIGN expr{...}|
2       identifier LPAREN call_args RPAREN{...}|
3       identifier{...}|
4       ADDR identifier{...}|
5       identifier LBRACKET expr RBRACKET {...}|
6       identifier LBRACKET expr RBRACKET ASSIGN expr{...}|

```

```

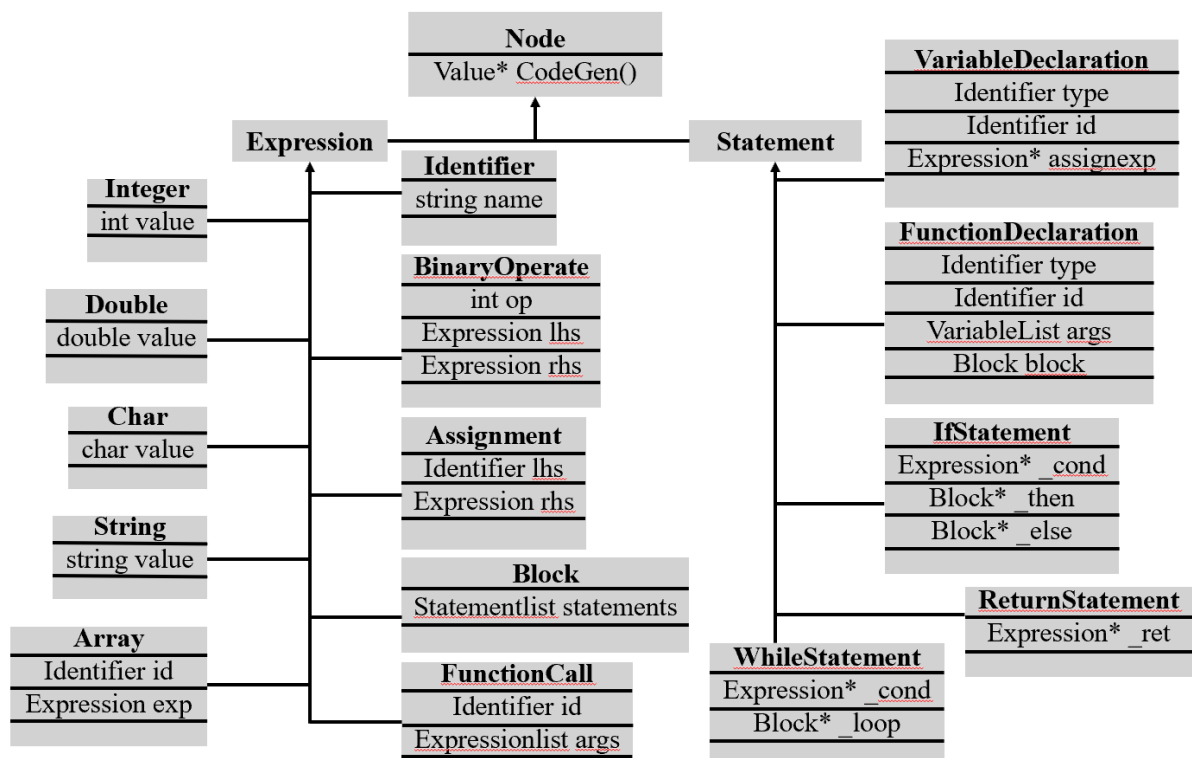
7      INTEGER {...} |
8      DOUBLE {...} |
9      CHARACTER {...} |
10     STRING {...} |
11     expr DIV expr{...} |
12     expr MUL expr{...} |
13     expr MOD expr{...} |
14     expr PLUS expr{...} |
15     expr MINUS expr{...} |
16     expr GT expr{...} |
17     expr GE expr{...} |
18     expr LT expr{...} |
19     expr LE expr{...} |
20     expr EQ expr{...} |
21     expr NEQ expr{...} |
22     expr AND expr{...} |
23     expr OR expr{...};

```

语义分析

AST

- 抽象语法树（Abstract Syntax Tree, AST）是源代码的抽象语法结构的树状表示。树上的每个节点都表示源代码中的一种结构，抽象语法树并不会表示出真实语法出现的每一个细节，比如嵌套括号被隐含在树的结构中，并没有以节点的形式呈现。
- 抽象语法树并不依赖于源语言的语法，也就是说语法分析阶段所采用的上下文无文文法，因为在写文法时，经常会对文法进行等价的转换（消除左递归，回溯，二义性等），这样会给文法分析引入一些多余的成分，对后续阶段造成不利影响，甚至会使合个阶段变得混乱。因此，很多编译器经常要独立地构造语法分析树，为前端，后端建立一个清晰的接口。
- Bison将Flex输入的token根据定义的文法及对应的action自动解析成AST的形式。
- 下面表示的AST为CMINUS Compiler的基本抽象语法树，由于篇幅问题省略了一些节点。



AST Node - Expression

- `Node` 为共同基类，所有的AST Node都继承自该类，该类有一个虚函数 `CodeGen()` 用于中间代码生成，不存在任何成员变量。

```
1 class Node {
2 public:
3     Node() {}
4     virtual ~Node() {}
5
6     virtual llvm::Value* CodeGen(CodeGenContext &ctx){
7         return nullptr;
8     }
9 };
```

- `Expression` 和 `Statement` 类为自 `Node` 类派生出的第一层AST节点，分别表示基本表达式和基本语句。表达式用于计算、语句用于执行操作，就像CPU的计算和控制功能。除此之外，由于大多数编程语言都是按照表达式和语句来设计的，所以将`expression`和`statement`来表示两个大类是合理的。
- `Integer`, `Char`, `Double`, `String` 继承自上一层节点 `Expression`，表示常量数据类型整型、字符型、双精度浮点型、字符串。每个节点会存储常量的具体的值，用于在代码生成时向global的堆中插入常量。
- `Array` 为数组，根据其标识符和大小调用`allocinst`进行内存分配。
- `Identifier` 的设计简洁地将所有的类型定义和标识符命名都统一到了一个子类里面。对于如 `int`, `char`, `double` 等的变量类型声明，`Identifier` 作为成员变量出现在其它的 `Expression` 类中。而对于 `a`, `x`, `y` 等的变量标识符，`Identifier` 会调用自身的 `CodeGen()` 将标识符插入符号表中，来进行变量和值的绑定。
- `BinaryOperate` 为基本的二元运算类，含有一个成员变量 `op` 表示运算表达式的操作类型，包括基本算术运算`+`, `-`, `*`, `/`, `%`，基本逻辑运算`&&`, `||`, `>=`, `>`, `<`, `<=`, `==`等。
- `Assignment` 为赋值运算类，包含一个表示左值的`Identifier`对象和一个表示右值的`Expression`对象，目标是把右值赋给左值。
- `Block` 为一个基本变量作用域代码块，它包含一个表示代码块中的语句的`StatementList`。其中`StatementList`的声明如下：

```
1 typedef std::vector<Statement *> StatementList;
```

- `FunctionCall` 为函数调用类，包含调用函数的名称和形式参数列表。其中`ExpressionList`的声明如下：

```
1 typedef std::vector<Expression *> ExpressionList;
```

AST Node - Statement

- `VariableDeclaration` 表示变量声明，包含一个 `Identifier` 类的类型和另一个 `Identifier` 类的变量名。同时为了同时兼容以下两种变量声明形式：

```
1 int x;
2 int x=0;
```

为该类增加成员变量 `Expression* assignexp` 表示在变量声明时是否直接对变量设置了初值。如果该指针为空，则未设置初值，反之。

- `FunctionDeclaration` 表示函数声明及定义。CMINUS Compiler目前仅支持了函数的声明和定义同时进行，不支持函数声明和定义的分离。在函数的声明和定义中，有 `Identifier` 类的成员变量 `type` 表示函数的返回类型，另一个 `Identifier` 类型的成员变量 `id` 表示函数名，以及 `VariableList` 类表示函数的形参列表类型，以及 `Block` 类的成员变量 `block` 表示函数体。

其中，`VariableList` 的声明如下：

```
1 typedef std::vector<VariableDeclaration*> VariableList;
```

- `IfStatement` 为条件语句。CMINUS Compiler目前仅支持如下格式的条件语句。

```
1 if (cond){
2     stmts
3 }
4 else{
5     stmts
6 }
```

条件语句由 `Expression` 类的 `_cond` 表示条件，`Block` 类的 `_then` 表示满足条件的代码块，`Block` 类的 `_else` 表示不满足条件的代码块。

- `WhileStatement` 为循环语句。CMINUS Compiler目前支持的while循环格式如下：

```
1 while(cond){
2     stmts
3 }
```

循环语句由 `Expression` 类的 `_cond` 表示条件，`Block` 类的 `_loop` 表示循环体的代码块。

- `ReturnStatement` 为返回语句。`Expression` 类型的指针 `_ret` 表示函数的返回值，若为空，表示无返回值。

中间代码生成

LLVM IR

LLVM IR有三种形式：

- 内存中的表示形式，如 `BasicBlock`，`Instruction` 这种cpp类。内存中IR模型其实就是对应LLVM实现中的class定义。
 - `Module`类，`Module`可以理解为一个完整的编译单元。一般来说，这个编译单元就是一个源文件，如一个后缀为cpp的源文件。
 - `Function`类，这个类顾名思义就是对应于一个函数单元。`Function`可以描述两种情况，分别是函数定义和函数声明。
 - `BasicBlock`类，这个类表示一个基本代码块，“基本代码块”就是一段没有控制流逻辑的基本流程，相当于程序流程图中的基本过程。
 - `Instruction`类，指令类就是LLVM中定义的基本操作，比如加减乘除这种算数指令、函数调用指令、跳转指令、返回指令等等。
- `bitcode`形式，这是一种序列化的二进制表示形式；
- LLVM汇编文件形式，这也是一种序列化的表示形式，与 `bitcode` 的区别是汇编文件是可读的、字符串的形式。

其中，LLVM汇编文件形式是CMINUS Compiler前端的主要目标。

AST To IR

对于AST定义的节点，需要调用LLVM的API根据其执行逻辑，生成对应的汇编IR表示。

在生成IR之前，首先要初始化一个IR在内存中的表示形式，即初始化一个IR生成的环境。

对于这个上下文环境，定义如下：

```
1  class CodeGenContext {
2
3      std::vector<CodeGenBlock *> symbol_stack; // symbol table stack
4
5  public:
6
7      llvm::Module *module; // store all functions and global variables
8      llvm::Function *cur_f; // indicate the current function
9      bool has_return; // indicate whether there is a return statement in the
current context
10     llvm::Value* ret_val=nullptr; //record the return value
11     std::vector<llvm::BasicBlock*> break_stack; // record the start place of
certain block containing break
12     std::string filename; // the target file to store the generated llvm IR
representation
13     llvm::BasicBlock* ret_bb; // return block of a function
14     llvm::Function* printf; // printf function prototype
15     llvm::Function* scanf; // scanf function prototype
16
17     CodeGenContext(std::string filename): filename(filename) {
18         module = new llvm::Module("main", global_ctx);
19         has_return = false;
20         printf = printf_prototype();
21         scanf = scanf_prototype();
22     }
23
24     void gen_code(Block *root); // function the generate IR representation
25
26     std::map<std::string, llvm::Value *> &get_sym_tab(); // function the get
the symbol table of the current context
27
28     void push_block(); // push a new environment to store symbol table
29
30     void pop_block(); // pop the top symbol table
31
32     llvm::Value *find_var(std::string var_name); // find variable in the
topest symbol table
33     llvm::Function* scanf_prototype(); // get prototype of scanf function
34     llvm::Function* printf_prototype(); // get prototype of printf function
35     llvm::BasicBlock* get_break(); // get break target basic block
36
37     void pop_break(); // pop basic block of break target
38
39     void insert_break(llvm::BasicBlock*); // insert target block of break
40
41 }
```

环境的成员变量如下

- `module: llvm::Module*`

Module代表了一块代码。它是一个比较完整独立的代码块，是一个最小的编译单元。每个Module含有函数,全局变量，符号表入口以及LLVM Linker等元素。

- `cur_f: llvm::Function*`

Function表示一个函数单元，可以描述两种情况，分别是函数定义和函数声明。`cur_f`表示当前函数，在编译时可以用来判断是否处在个函数的定义内部。

- `ret_val: llvm::Value*`

用于记录函数的返回值。

- `break_stack: std::vector<llvm::BasicBlock*>`

用于定位break语句分支跳转的目标地址。

Code Generation

- Integer, Double, Char等基本常量类型

直接返回Constant类型指令，即在module的开头就定义这个全局变量。

```
1 return ConstantInt::get(Type::getInt32Ty(global_ctx), value, true);
2 return ConstantFP::get(Type::getDoubleTy(global_ctx), value);
3 return builder.getInt8(value);
```

- Identifier 变量

首先需要判断该变量是否定义。若在当前symbol table中找到了该变量，则已经定义，否则抛出异常。

```
1 // find the variable in the symbol table of current context
2 value *var = ctx.find_var(name);
3 if (var == nullptr) {
4     throw std::logic_error("undeclared variable!" + name + "\n");
5 }
```

根据符号表中存储的变量类型来load这个变量的值。

```
1 res = new LoadInst(type, var, "LoadInst", false,
2 builder.GetInsertBlock());
```

- BinaryOperate 算数运算

首先需要调用Expression的 `CodeGen()` 得到左值和右值的值类型。

```
1 value* left = lhs.CodeGen(ctx);
2 value* right = rhs.CodeGen(ctx);
```

根据左值和右值的类型进行适当的类型转换，目前CMINUS Compiler仅支持了char2int的类型转换。通过调用 `CreateCast` 和 `ZExt` 指令来实现类型转换。

```

1  Instruction::CastOps type_inst(Type* src, Type* dst){
2      if(src==Type::getInt8Ty(global_ctx) &&
   dst==Type::getInt32Ty(global_ctx))
3          return Instruction::ZExt;
4      throw std::logic_error("[ERROR] Wrong typecast");
5  }
6
7  Value* type_convert(Value* src, Type* dst){
8      Instruction::CastOps op = type_inst(src->getType(), dst);
9      return builder.CreateCast(op, src, dst, "typeconvert");
10 }

```

经过类型检查后，需要判断左值和右值是否为浮点型。若为浮点型需要调用浮点数运算指令，否则为整型运算。下以PLUS运算为例：

```

1  bool is_double = left->getType()->isDoubleTy();
2  case PLUS:
3      if (is_double)
4          return builder.CreateFAdd(left, right);
5      else
6          return builder.CreateAdd(left, right);

```

- Assignment 赋值运算

首先需要在symbol table中寻找左边的identifier

```

1  // find the identifier in the symbol stack
2  Value* res = ctx.find_var(lhs.name);

```

然后调用Expression的 `CodeGen()` 生成其值，并使用Store指令进行存储。

```

1  Value* right=rhs.CodeGen(ctx);
2  if (right->getType() != res->getType()->getPointerElementType())
3      throw std::logic_error("assignment type not match!\n");
4
5  return builder.CreateStore(right, res);

```

- Block 代码块代码生成

对于处于block内的statements逐个调用其 `CodeGen()` 函数，

```

1  for(; ite!=statements.end(); ite++){
2      res=(*ite)->CodeGen(ctx);
3      if (ctx.has_return == true)
4          break;
5  }

```

- FunctionCall 函数调用

首先判断函数是否是 `scanf` 和 `printf`，如果是则调用其prototype函数原型，否则从环境中寻找Function进行调用。

```

1  if (id.name=="printf"){
2      return printf_gen(ctx, args);
3  }
4  else if (id.name=="scanf"){
5      return scanf_gen(ctx, args);
6  }
7  // check function in the context
8      Function* func = ctx.module->getFunction(id.name.c_str());

```

计算参数列表的值

```

1  // calculate arg value of the passing parameters
2      std::vector<Value*> args_list;
3      for (auto ite : args)
4          args_list.push_back((*ite).CodeGen(ctx));

```

调用函数

```

1  CallInst* call = CallInst::Create(func, makeArrayRef(args_list), id.name,
    builder.GetInsertBlock());

```

- **VariableDeclaration** 变量声明

首先根据 `cur_f` 判断该变量是否为全局变量，如果是则初始化一个全局变量

```

1  GlobalVariable* var = new GlobalVariable(*(ctx.module), llvm_type, false,
    \      GlobalValue::PrivateLinkage, 0, id.name);

```

若不是，则将变量计入当前环境的栈顶符号表，并为变量分配栈空间。

```

1  AllocaInst *inst = new AllocaInst(llvm_type, blk->getParent()-
    >getParent()->getDataLayout().getAllocaAddrSpace(), id.name.c_str(),
    blk);
2
3  // insert the newly defined variables into the symbol table
4  ctx.get_sym_tab()[id.name] = inst;

```

判断在变量声明时是否对变量赋初值，若赋初值则需调用 `Assignment` 节点的 `CodeGen()` 函数。

```

1  if(assignexp!=NULL){
2      Assignment asgn(id, *assignexp);
3      asgn.CodeGen(ctx);
4  }

```

- **FunctionDeclaration** 函数声明

向环境的Module中加入这个函数。

```

1  Function* func = Function::Create(func_type,
    GlobalValue::ExternalLinkage, id.name.c_str(), ctx.module);

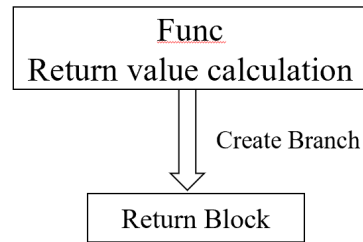
```

首先初始化两个Basic Block，分别为函数体BasicBlock和Return语句的BasicBlock。

```

1 BasicBlock* func_bb = BasicBlock::Create(global_ctx, "entry", func, 0);
2 BasicBlock* ret_bb = BasicBlock::Create(global_ctx, "return", func, 0);

```



压栈一个新的symbol table以新建一个函数体的变量空间，处理函数形参列表的赋值，并对函数体的代码调用 `Block` 节点的 `CodeGen()` 函数进行IR生成。

```

1 ctx.push_block();
2 ctx.cur_f = func;
3 builder.SetInsertPoint(func_bb);
4 for(auto ite : args){
5     (*ite).CodeGen(ctx);
6     arg_val = &*(arg_ite++);
7     arg_val->setName(ite->id.name.c_str());
8     StoreInst* inst = new StoreInst(arg_val, ctx.get_sym_tab()[ite->id.name], false, func_bb);
9 }
10 block.CodeGen(ctx);

```

处理函数返回，将栈顶symbol table弹出，恢复原来的变量作用域，计算返回值，并创建ret指令

```

1 if(type.name == "void") {
2     builder.CreateRetVoid();
3 } else {
4     if(ctx.ret_val == nullptr){
5         throw std::logic_error("return is needed for non void function!\n");
6     }
7     Value* ret = builder.CreateLoad(get_type(type.name, false), ctx.ret_val, "");
8     builder.CreateRet(ret);
9 }
10 ctx.pop_block();
11 ctx.cur_f = nullptr;

```

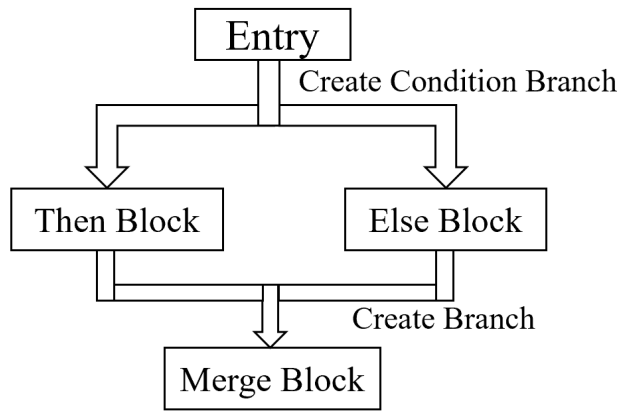
- `IfStatement` 条件语句

初始化三个BasicBlock: then_bb, else_bb, merge_bb。

```

1 BasicBlock* then_bb = BasicBlock::Create(global_ctx, "if", func);
2 BasicBlock* else_bb = BasicBlock::Create(global_ctx, "else", func);
3 BasicBlock* merge_bb = BasicBlock::Create(global_ctx, "merge", func);

```



根据_cond创建条件跳转。

```

1 Value* cond = _cond->CodeGen(ctx);
2 cond = builder.CreateICmpNE(cond,
  ConstantInt::get(Type::getInt1Ty(global_ctx), 0, true), "ifcond");

```

将builder指向then block的入口，新建一个变量作用域，处理_then语句块。

```

1 builder.SetInsertPoint(then_bb);
2 ctx.push_block();
3 _then->CodeGen(ctx);
4 ctx.pop_block();
5 if(ctx.has_return)
6     ctx.has_return=false;
7 else
8     builder.CreateBr(merge_bb);

```

将builder指向else block的入口，新建一个变量作用域，处理_else语句块。

```

1 builder.SetInsertPoint(else_bb);
2 if (_else!=nullptr){
3     ctx.push_block();
4     _else->CodeGen(ctx);
5     ctx.pop_block();
6
7     if (ctx.has_return)
8         ctx.has_return=false;
9     else
10        builder.CreateBr(merge_bb);

```

在两个block的末尾都要创建分支指令，跳到merge block。

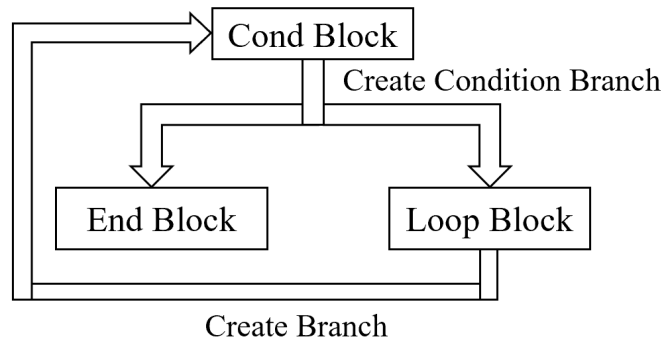
- whileStatement 循环语句

初始化三个BasicBlock: cond_bb, loop_bb, end_bb。

```

1 BasicBlock* cond_bb = BasicBlock::Create(global_ctx, "while_cond", func);
2 BasicBlock* loop_bb = BasicBlock::Create(global_ctx, "while_loop", func);
3 BasicBlock* end_bb = BasicBlock::Create(global_ctx, "while_end", func);

```



处理_cond, 加入条件跳转指令, 根据condition分别跳转至end block和loop block。

```

1 value* cond = _cond->CodeGen(ctx);
2 cond = builder.CreateICmpNE(cond,
  ConstantInt::get(Type::getInt1Ty(global_ctx), 0, true), "while_cond");
3 auto branch = builder.CreateCondBr(cond, loop_bb, end_bb);

```

处理_loop, 首先调用Block节点的CodeGen()函数生成指令, 再创建branch指令跳转至cond block重新进行条件判断。

```

1 builder.SetInsertPoint(loop_bb);
2 ctx.push_block();
3 _loop->CodeGen(ctx);
4 ctx.pop_block();
5 ctx.pop_break();

```

- ReturnStatement 返回语句

根据是否有返回值调用 Expression 节点的``CodeGen()``函数, 再无条件跳转到函数的ret block。

```

1 if (_ret==nullptr){
2     if(ctx.cur_f->getReturnType()->isVoidTy())
3         ;
4     else
5         throw std::logic_error("expected return value for non void
function!");
6 }else {
7     std::cout << "return expression: " << typeid(*_ret).name()<<
std::endl;
8     value *ret = _ret->CodeGen(ctx);
9     if (ret->getType() != ctx.cur_f->getReturnType())
10        throw std::logic_error("type of return value must match the
function type!");
11    builder.CreateStore(ret, ctx.ret_val);
12 }
13 ctx.has_return=true;
14 return builder.CreateBr(ctx.ret_bb);

```

- Address 取地址

直接在符号表中查找并返回identifier对应的Value*。

```

1  std::cout << "Get Address of identifier "<<id.name<<std::endl;
2  value* res = ctx.find_var(id.name);
3  if(res==nullptr)
4      throw std::logic_error("undeclared variable "+id.name+" \n");
5  else
6      return res;

```

- **Array** 数组元素访问

主要是需要知道identifier的地址，及元素索引，来建立Load指令进行访问。

```

1  value* ptr = builder.CreateInBoundsGEP(val->getType()-
    >getPointerElementType(), val, \
2                                     ArrayRef<Value*>(idx_list),
    "tmparray");
3  return builder.CreateLoad(ptr->getType()->getPointerElementType(), ptr,
    "tmpvar");

```

- **ArrayAssignment** 数组元素赋值

类似 Assignment 节点，只是左值为数组元素，调用数组元素访问进行store。

```

1  value* lhs = builder.CreateInBoundsGEP(val->getType()-
    >getPointerElementType(), val,
2                                     ArrayRef<Value*>(idx_list),
    "tmparray");
3  value* rhs = _exp.CodeGen(ctx);
4  builder.CreateStore(rhs, lhs);

```

测试

Make Process

```

1  cmake -S ./build
2  cd build
3  make
4  ./Compiler < test.cpp
5  ./run.sh

```

Basic Test

- test_prime
 - source code

```

1  int main(){
2      int i;
3      int a=2;
4      int tmp;
5      scanf("%d", &i);
6
7      if(i==1){
8          printf("0");
9          return 0;
10     }

```



```

11     else{
12         tmp=0;
13     }
14
15     if(i==2){
16         printf("1");
17         return 0;
18     }
19     else{
20         tmp=0;
21     }
22
23     while(a<i/2+1){
24         if (i%a==0){
25             printf("0");
26             return 0;
27         }
28         else{
29             tmp=0;
30         }
31         a=a+1;
32     }
33     printf("1");
34     return 0;
35 }

```

- o IR representation

```

1  ; ModuleID = 'main'
2  source_filename = "main"
3
4  @_Const_String_ = private constant [3 x i8] c"%d\00"
5  @_Const_String_.1 = private constant [2 x i8] c"0\00"
6  @_Const_String_.2 = private constant [2 x i8] c"1\00"
7  @_Const_String_.3 = private constant [2 x i8] c"0\00"
8  @_Const_String_.4 = private constant [2 x i8] c"1\00"
9
10 declare i32 @printf(i8*, ...)
11
12 declare i32 @scanf(...)
13
14 define i32 @main() {
15 entry:
16     %0 = alloca i32, align 4
17     %i = alloca i32, align 4
18     %a = alloca i32, align 4
19     store i32 2, i32* %a, align 4
20     %tmp = alloca i32, align 4
21     %scanf = call i32 (...) @scanf(i8* getelementptr inbounds ([3 x
22 i8], [3 x i8]* @_Const_String_, i32 0, i32 0), i32* %i)
23     %LoadInst = load i32, i32* %i, align 4
24     %l = icmp eq i32 %LoadInst, 1
25     %ifcond = icmp ne i1 %l, false
26     br i1 %ifcond, label %if, label %else
27
28 return:                                ; preds =
29     %while_end, %if10, %if1, %if

```

```

28     %2 = load i32, i32* %0, align 4
29     ret i32 %2
30
31 if:                                     ; preds = %entry
32     %printf = call i32 @printf(i8*, ...) @printf(i8* getelementptr inbounds
33 ([2 x i8], [2 x i8]* @_Const_String_.1, i32 0, i32 0))
34     store i32 0, i32* %0, align 4
35     br label %return
36
37 else:                                   ; preds = %entry
38     store i32 0, i32* %tmp, align 4
39     br label %merge
40
41 merge:                                 ; preds = %else
42     %LoadInst4 = load i32, i32* %i, align 4
43     %3 = icmp eq i32 %LoadInst4, 2
44     %ifcond5 = icmp ne i1 %3, false
45     br i1 %ifcond5, label %if1, label %else2
46
47 if1:                                   ; preds = %merge
48     %printf6 = call i32 @printf(i8*, ...) @printf(i8* getelementptr inbounds
49 ([2 x i8], [2 x i8]* @_Const_String_.2, i32 0, i32 0))
50     store i32 0, i32* %0, align 4
51     br label %return
52
53 else2:                                 ; preds = %merge
54     store i32 0, i32* %tmp, align 4
55     br label %merge3
56
57 merge3:                                ; preds = %else2
58     br label %while_cond
59
60 while_cond:                            ; preds =
61     %merge12, %merge3
62     %LoadInst7 = load i32, i32* %a, align 4
63     %LoadInst8 = load i32, i32* %i, align 4
64     %4 = sdiv i32 %LoadInst8, 2
65     %5 = add i32 %4, 1
66     %6 = icmp slt i32 %LoadInst7, %5
67     %while_cond9 = icmp ne i1 %6, false
68     br i1 %while_cond9, label %while_loop, label %while_end
69
70 while_loop:                            ; preds =
71     %while_cond
72     %LoadInst13 = load i32, i32* %i, align 4
73     %LoadInst14 = load i32, i32* %a, align 4
74     %7 = srem i32 %LoadInst13, %LoadInst14
75     %8 = icmp eq i32 %7, 0
76     %ifcond15 = icmp ne i1 %8, false
77     br i1 %ifcond15, label %if10, label %else11
78
79 while_end:                             ; preds =
80     %while_cond
81     %printf18 = call i32 @printf(i8*, ...) @printf(i8* getelementptr inbounds
82 ([2 x i8], [2 x i8]* @_Const_String_.4, i32 0, i32 0))
83     store i32 0, i32* %0, align 4
84     br label %return
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99

```

```

80 if10:                                     ; preds =
    %while_loop
81     %printf16 = call i32 @printf(i8* getelementptr inbounds
    ([2 x i8], [2 x i8]* @_Const_String_.3, i32 0, i32 0))
82     store i32 0, i32* %, align 4
83     br label %return
84
85 else11:                                    ; preds =
    %while_loop
86     store i32 0, i32* %tmp, align 4
87     br label %merge12
88
89 merge12:                                  ; preds = %else11
90     %LoadInst17 = load i32, i32* %, align 4
91     %9 = add i32 %LoadInst17, 1
92     store i32 %9, i32* %, align 4
93     br label %while_cond
94 }

```

- test result

```

wjj@LAPTOP-F02JHQPE:/mnt/e/compiler_wjj/build$ ../../Compiler/final_proj/bin/linux/amd64/test_prime ./test_exe/test_prime
Running test cases for ./test_exe/test_prime
Sample 1      PASSED
Sample 2      PASSED
Sample 3      PASSED
Test 4        PASSED
Test 5        PASSED
Passed        5/5
2023-05-29 02:38:51

```

- test_hanoi
 - source code

```

1  int move(char x, char y)
2  {
3      int c=0;
4      printf("%c->%c\n",x,y);
5      return c;
6  }
7
8  int hanoi(int n, char one, char two, char three)
9  {
10     int tmp=1;
11     if(n==tmp) {
12         move(one, three);
13     }
14     else
15     {
16         hanoi(n-tmp, one, three, two);
17         move(one, three);
18         hanoi(n-tmp, two, one, three);
19     }
20     return 0;
21 }
22
23 int main()
24 {

```

```

25     int m;
26     int d=0;
27     scanf("%d",&m);
28     char a='A';
29     char b='B';
30     char c='C';
31     hanoi(m,a,b,c);
32     return 0;
33 }
34

```

- test result

```

wjj@LAPTOP-F02JHQPE:/mnt/e/compiler_wjj/build$ ../../Compiler/final_proj/bin/linux/amd64/test_hanoi ./test_exe/test_hano
Running test cases for ./test_exe/test_hano
Sample 1          PASSED
Sample 2          PASSED
Test 3            PASSED
Test 4            PASSED
Passed            4/4
2023-05-29 02:39:58

```

- test_matrix

- source code

```

1  int main()
2  {
3      int M;
4      int N;
5      int P;
6      int a[100];
7      int b[100];
8      int i=0;
9      int j=0;
10     int k=0;
11     int sum=0;
12     int tmp1;
13     int tmp2;
14     int tmp;
15     scanf("%d %d %d", &M, &N, &P);
16     while(i<M*N){
17         scanf("%d", &tmp);
18         a[i]=tmp;
19         i = i+1;
20     }
21
22     i=0;
23     while(i<N*P){
24         scanf("%d", &tmp);
25         b[i]=tmp;
26         i = i+1;
27     }
28     i=0;
29     while(i<M){
30         j=0;
31         while(j<P){
32             sum=0;
33             k=0;
34             while(k<N){

```

```

35         tmp1 = a[N*i+k];
36         tmp2 = b[P*k+j];
37         sum = sum + tmp1*tmp2;
38         k = k+1;
39     }
40     if(j==0) {
41         printf("%d", sum);
42     }
43     else {
44         printf(" %d", sum);
45     }
46     j = j+1;
47 }
48 printf("\n");
49 i = i+1;
50 }
51
52 return 0;
53 }

```

- test result

```

wj@LAPTOP-F02JHQPE:/mnt/e/compiler_wj/build$ ../../Compiler/final_proj/bin/linux/amd64/test_matrix ./test_exe/test_matrix
Running test cases for ./test_exe/test_matrix
Sample 1      PASSED
Test 2       PASSED
Test 3       PASSED
Test 4       PASSED
Passed       4/4
2023-05-29 02:41:28

```

- test_lcs

- source code

```

1
2 char a[20];
3 char b[20];
4 int c[400];
5 int n=0;
6 int m=0;
7 int i=0;
8 int j=0;
9 int k=0;
10 int tmp1;
11 int tmp2;
12
13 int main()
14 {
15
16     char tmp=a[19];
17     scanf("%s %s",a,b);
18
19     while(a[i]>tmp){
20         i=i+1;
21         n=n+1;
22     }
23     i=0;
24     while(b[i]>tmp){
25         i=i+1;

```

```

26         m=m+1;
27     }
28     i=0;
29     while(i<n+1){
30         tmp1=m+1;
31         tmp1=i*tmp1;
32         c[tmp1]=0;
33         i=i+1;
34     }
35     i=0;
36     while(i<m+1){
37         c[i]=0;
38         i=i+1;
39     }
40     i=1;
41     while(i<n+1){
42         j=1;
43         while(j<m+1){
44             if(a[i-1]==b[j-1])
45             {
46                 tmp1=m+1;
47                 tmp1=i*tmp1;
48                 tmp1=tmp1+j;
49                 tmp2=m+1;
50                 k=i-1;
51                 tmp2=tmp2*k;
52                 tmp2=tmp2+j-1;
53                 c[tmp1]=c[tmp2]+1;
54             }
55             else
56             {
57                 tmp1=i-1;
58                 k=m+1;
59                 tmp1=k*tmp1;
60                 tmp1=tmp1+j;
61                 tmp1=c[tmp1];
62                 tmp2=m+1;
63                 tmp2=tmp2*i;
64                 tmp2=tmp2+j;
65                 tmp2=tmp2-1;
66                 tmp2=c[tmp2];
67                 k=m+1;
68                 k=k*i;
69                 k=k+j;
70                 if(tmp1>tmp2){
71                     c[k]=tmp1;
72                 }
73                 else{
74                     c[k]=tmp2;
75                 }
76             }
77             j=j+1;
78         }
79         i=i+1;
80     }
81     tmp1=m+1;
82     tmp1=tmp1*n;
83     tmp1=tmp1+m;

```

```

84     printf("%d\n",c[tmp1]);
85
86     return 0;
87 }

```

- test result

```

wjj@LAPTOP-F02JHQPE:/mnt/e/compiler_wjj/build$ ../../Compiler/final_proj/bin/linux/amd64/test_lcs ./test_exe/test_maxseq
Running test cases for ./test_exe/test_maxseq
Sample 1      PASSED
Sample 2      PASSED
Test 3        PASSED
Test 4        PASSED
Test 5        PASSED
Test 6        PASSED
Test 7        PASSED
Test 8        PASSED
Passed        8/8
2023-05-29 02:42:09

```

- test_arithmetic

- source code

```

1  int num;
2  int tmp;
3  int num_cnt=0;
4  int sym_cnt=0;
5  int num_stack[15];
6  int symbol_stack[15];
7  int num_tmp[10];
8  int cnt=0;
9  int p=2;
10 int q;
11 char cc[100];
12
13
14 int calc(){
15     int i;
16     int res;
17     int tmp2;
18     i=symbol_stack[sym_cnt-1];
19     tmp=num_stack[num_cnt-2];
20     tmp2=num_stack[num_cnt-1];
21     if(i==43){
22         res=tmp+tmp2;
23     }
24     else{
25         if(i==45){
26             res=tmp-tmp2;
27         }
28         else{
29             if(i==42){
30                 res=tmp*tmp2;
31             }
32             else{
33                 if(i==47){
34                     res=tmp/tmp2;
35                 }
36                 else{
37                     res=0;

```

```

38         }
39     }
40 }
41 }
42 sym_cnt=sym_cnt-1;
43 num_cnt=num_cnt-1;
44 num_stack[num_cnt-1]=res;
45
46 return 0;
47 }
48
49 int calc_pri(){
50     int i;
51     if(sym_cnt==0){
52         p=2;
53         if(q==1){
54             p=0;
55         }
56         else{
57             tmp=0;
58         }
59     }
60     else{
61         i=symbol_stack[sym_cnt-1];
62         if(i==42){
63             p=4;
64         }
65         else{
66             if(i==45){
67                 p=4;
68             }
69             else{
70                 p=3;
71             }
72         }
73     }
74     return 0;
75 }
76
77 int deal_sym(){
78     if(num==42) {
79         q=4;
80     }
81     else {
82         if(num==47) {
83             q = 4;
84         }
85         else {
86             if (num == 10) {
87                 q = 1;
88             } else {
89                 q = 3;
90             }
91         }
92     }
93
94     while(p+1>q){
95         calc();

```



```

96     calc_pri();
97 }
98 symbol_stack[sym_cnt]=num;
99 sym_cnt=sym_cnt+1;
100 p=q;
101
102     return 0;
103 }
104
105 int calc_num(){
106     int res=0;
107     int i=0;
108     while(i<cnt){
109         res=10*res;
110         res=res+num_tmp[i];
111         i=i+1;
112     }
113     cnt=0;
114     return res;
115 }
116 int main(){
117     char a=cc[0];
118     int yy[100];
119     int l;
120     int j=0;
121     scanf("%s", cc);
122     tmp=0;
123     while(cc[tmp]>a){
124         yy[tmp]=cc[tmp]-0;
125         tmp=tmp+1;
126     }
127     tmp=tmp-1;
128     if(yy[tmp]==10){
129         num=0;
130     }
131     else{
132         tmp=tmp+1;
133         yy[tmp]=10;
134     }
135     l=tmp+1;
136     while(j<l) {
137         num=yy[j];
138         if(num>47){
139             num_tmp[cnt]=num-48;
140             cnt=cnt+1;
141         }
142         else{
143             tmp = calc_num();
144             num_stack[num_cnt] = tmp;
145             num_cnt = num_cnt + 1;
146             deal_sym();
147         }
148
149         if(num==10){
150             printf("%d\n", num_stack[0]);
151             return 0;
152         }
153         else{

```

```

154         tmp=0;
155     }
156     j=j+1;
157 }
158 return 0;
159 }

```

- test result

```

wj@LAPTOP-F02JHQPE:/mnt/e/compiler_wjj/build$ ../../Compiler/final_proj/bin/linux/amd64/test_arithmetic ./test_exe/test_arithmetic
Running test cases for ./test_exe/test_arithmetic
Sample 1      PASSED
Sample 2      PASSED
Sample 3      PASSED
Test 4        PASSED
Test 5        PASSED
Test 6        PASSED
Test 7        PASSED
Test 8        PASSED
Passed       8/8
2023-05-29 02:43:17

```

Error Test

- Lexer Error
 - unknown token

```

1  int main(){
2      char ^a;
3      int c=1;
4      scanf("%c",&a);
5      return 0;
6  }

```

```

terminate called after throwing an instance of 'std::logic_error'
  what():  Unknown token!

Aborted

```

- Syntax Error
 - BRACE not match, need `{ }` in if _then block

```

1  int main(){
2      int a=1;
3      if (a==1)
4          printf("%d", a);
5      return 0;
6  }

```

```

ERROR: syntax error, unexpected IDENTIFIER, expecting LBRACE

```

- PAREN not match

```

1  int main(){
2      int a=1;
3      while a>0){
4          a=a-1;
5      }
6      return 0;
7  }

```

ERROR: syntax error, unexpected IDENTIFIER, expecting LPAREN

- Function Decalaration arg list error

```

1  int f int i){
2      i=1;
3      return i;
4  }
5  int main(){
6      int i=0;
7      printf("%d",f(i));
8      return 0
9  }

```

ERROR: syntax error, unexpected IDENTIFIER, expecting LPAREN or LBRACKET or SEMI or ASSIGN

- Semantic Error

- type not defined

```

1  int main(){
2      float a=1.0;
3      return 0;
4  }

```

```

terminate called after throwing an instance of 'std::logic_error'
  what():  Type float not defined!

Aborted

```

- variable not defined

```

1  int main(){
2      if(i>0){
3          i=1;
4      }
5      else{
6          i=2;
7      }
8      return 0;
9  }

```

```

terminate called after throwing an instance of 'std::logic_error'
  what():  undeclared variable!i

```

- function arg type not match

```

1  int f(int m, int n){
2      m=1;
3      printf("%d\n",m);
4      return 0;
5  }
6
7  int main(){
8      int a;
9      scanf("%c",&a);
10     f(a);
11
12     return 0;
13 }

```

```

terminate called after throwing an instance of 'std::logic_error'
  what():  function arg type not match

Aborted

```

- function not defined

```

1  int main(){
2      char a;
3      scanf("%c",&a);
4      f(a);
5
6      return 0;
7  }

```

```

terminate called after throwing an instance of 'std::logic_error'
  what():  cannot find function: f

```

- calculation type not match

```

1  int main(){
2      char a;
3      int c=1;
4      scanf("%c",&a);
5      if(a>c){
6          printf("in if\n");
7      }
8      else{
9          printf("in else\n");
10     }
11     return 0;
12 }

```

```

terminate called after throwing an instance of 'std::logic_error'
  what():  [ERROR] Wrong typecast

```

- local variable redefined

```

1  int main(){
2      char a;
3      int c=1;
4      scanf("%c",&a);
5      printf("%c\n",a);
6      int a=1;
7      printf("%d\n",a);
8
9      return 0;
10 }

```

```

terminate called after throwing an instance of 'std::logic_error'
  what():  Local variable has been defined! a
Aborted

```

- no return value for non-void function

```

1  int g(char x){
2      printf("%c\n",x);
3      return;
4  }
5  int main(){
6      char a;
7      int c=1;
8      scanf("%c",&a);
9      g(a);
10
11     return 0;
12 }

```

```

terminate called after throwing an instance of 'std::logic_error'
  what():  expected return value for non void function!
Aborted

```

- return value type not match for function

```

1  int g(char x){
2      printf("%c\n",x);
3      return 1.0;
4  }
5  int main(){
6      char a;
7      int c=1;
8      scanf("%c",&a);
9      g(a);
10
11     return 0;
12 }

```

```

terminate called after throwing an instance of 'std::logic_error'
  what():  type of return value must match the function type!
Aborted

```

- o illegal use of `break`

```
1 int main(){
2     char a;
3     int c=1;
4     scanf("%c",&a);
5     if(a=='a'){
6         break;
7     }
8     return 0;
9 }
```

```
terminate called after throwing an instance of 'std::logic_error'
  what():  break should be used in loops!
Aborted
```

问题与挑战

- `scanf`和`printf`的函数调用
 - o 由于`scanf`和`printf`是c中的库函数，CMINUS Compiler目前不支持 `#include<>` 形式的C本身代码库的链接，目前仅通过调用 `llvm` 内置的C函数外部链接函数接口实现。

```
1 llvm::Function *func = llvm::Function::Create(
2     scanf_type, llvm::Function::ExternalLinkage,
3     llvm::Twine("scanf"),
4     module
5 );
6 func->setCallingConv(llvm::CallingConv::C);
```

具体的参考了github中llvm的[demo](#)

- o 当我们直接把以上函数原型封装成函数获取时，会发现当程序中调用多次 `printf` 或 `scanf` 会发生找不到函数的错误。

```
1 int main(){
2     char c1;
3     char c2;
4     scanf("%c", &c1);
5     scanf("%c", &c2);
6     return 0;
7 }
```

观察生成的IR，会发现由于调用 `prototype` 获取函数的时候，创建函数指定的名字为 `scanf`，都偶此调用，llvm默认为重名函数，会自动把函数名字按照 `scanf`，`scanf1` 重新命名，导致函数调用错误。由此，可以发现，`printf` 和 `scanf` 在同一个module中只能创建一次，因此应当被作为上下文环境的成员变量直接存放。

- BasicBlock的执行逻辑
 - o 在条件语句、循环语句、函数调用的执行过程中，一个代码块无可避免的会有多个 terminator，需要合理的划分BasicBlock，使得这些语句拥有正常的执行逻辑。
- 当定义多个变量的时候，CMINUS Compiler会报 `Segementation Fault`。但尝试将局部变量改成全局变量，允许定义的变量数目上限增加，此问题大概是由于程序运行被分配到的Memory大小的问题，但一个合格的编译器应该允许定义很多变量，目前不知道有什么解决方法。

- 由于时间有限，CMINUS Compiler在设计的时候还存在很多不完善的地方。比如还未调试条件语句不匹配else的情形，while中使用break，int2double, char2int, char2double等基本类型之间的类型转换，以及代码优化，去掉不必要的执行逻辑，一维数组扩展成多维等。