

Stage 6 报告

程思翔 2021010761

写在前面:

这一个 Stage 的实现中对着前中后端的报错疯狂调 Bug, 基本实现过程就是 [构造测例] -> [看报错] -> [调代码], 好在终于是实现地较为完备, 痛并快乐着 (x).

报告或许会较长, 尽管尝试简化对实现思路的叙述, 最后还是保留了如下内容, 以便使用尽可能简洁的描述, 将实现过程清晰地展示出来.

1 实验内容

1.1 step 10

词法语法分析: 修改 `frontend/ast/tree.py` 中 `Program` 节点的定义, 并添加 `globalVars()` 方法传递全局变量键值对. 在 `frontend/parser/ply_parser.py` 中参考已有实现, 给出相关 CFG 文法:

```
1 # frontend/ast/tree.py
2 class Program(ListNode[Union["Function", "Declaration"]]):
3     ...
4     def globalVars(self) -> dict[str, int]:
5         return {decl.getattr('symbol').name:
6                 decl.getattr('symbol').initValue for decl in self if
7                 isinstance(decl, Declaration)}
```

语义分析: 只需修改 `frontend/typecheck/namer.py` 的 `visitDeclaration` 部分, 如果当前作用域为全局作用域, 修改 `symbol.isGlobal = True` 即可, 并为全局变量设置初始值 `symbol.initValue`.

中间代码生成: 在 `utils/tac/tacinstr.py` 中参照已有实现, 添加全局变量地址加载、全局变量加载和存储的 TAC 指令:

```

1  # utils/tac/tacinstr.py
2  class LoadAddress(TACInstr):
3      # LOAD SYMBOL ADDRESS...
4  class LoadIntLiteral(TACInstr):
5      # LOAD SYMBOL...
6  class StoreIntLiteral(TACInstr):
7      # STORE SYMBOL...

```

在 `utils/tac/tacgen.py` 中 `TACFuncEmitter` 类实现与之有关的 `Visitor` 模式方法, 修改 `TACGen` 类的 `visitAssignment` 及 `visitIdentifier` 方法, 判断访问标识符和赋值操作是否对全局变量进行:

```

1  # utils/tac/tacgen.py
2  class TACGen(Visitor[TACFuncEmitter, None]):
3      def visitIdentifier(self, ident: Identifier, mv:
TACFuncEmitter) -> None:
4          #! 标识符是全局变量
5          if ident.getattr("symbol").isGlobal:
6              ident.setattr('val',
mv.visitLoadIntLiteral(ident.getattr('symbol')))
7          #! 否则
8          else:
9              ident.setattr('val', ident.getattr('symbol').temp)
10
11      def visitAssignment(self, expr: Assignment, mv:
TACFuncEmitter) -> None:
12          expr.rhs.accept(self, mv)
13          #! 左值是全局变量
14          if expr.lhs.getattr('symbol').isGlobal:
15              mv.visitStoreIntLiteral(expr.lhs.getattr('symbol'),
expr.rhs.getattr("val"))
16          #! 否则
17          else:
18              expr.lhs.accept(self, mv)
19              expr.setattr(

```

```

20         "val",
mv.visitAssignment(expr.lhs.getattr("symbol").temp,
expr.rhs.getattr("val"))
21     )

```

目标代码生成: 在 `utils/riscv.py` 中参照已有实现, 添加全局变量地址获取、加载和存储的 RISC-V 指令:

```

1  # utils/riscv.py
2  class LoadAddress(TACInstr):
3      # LOAD SYMBOL ADDRESS...
4  class LoadIntLiteral(TACInstr):
5      # LOAD SYMBOL...
6  class StoreIntLiteral(TACInstr):
7      # STORE SYMBOL...

```

在 `backend/riscv/riscvasmemitter.py` 中, `RiscvAsmEmitter` 对象初始化时打印全局变量至 `.data` 区, 并为 `RiscvInstrSelector` 实现相应全局变量指令访问方法:

1.2 step 11

词法语法分析: 修改 `frontend/ast/tree.py`, 通过 `Program` 的 `globalVars()` 方法返回全局变量与全局数组键值对; 修改 `Declaration` 节点定义, 添加数组维度声明 `init_dim`; 同时添加索引运算节点 `IndexExpr`:

```

1  # frontend/ast/tree.py
2  class IndexExpr(Expression):
3      def __init__(self, base: Expression, index: Expression) ->
None:
4      super().__init__("index_expr")
5      self.base = base
6      self.index = index

```

在 `frontend/parser/ply_parser.py` 中参考已有实现, 给出相关 CFG 文法.

语义分析: 修改 `frontend/typecheck/namer.py`, `visitFunction` 中记录当前函数声明的局部数组; 在 `visitDeclaration` 中依据标识符类型进行初始化, 如果是数组则调用 `ArrayType.multidim`; 添加 `visitIndexExpr` 中递归实现数组索引访问.

```
1  # frontend/typecheck/namer.py
2  def visitDeclaration(self, decl: Declaration, ctx: ScopeStack) ->
    None:
3      ...
4      if decl.init_dim:
5          decl_type = ArrayType.multidim(decl.var_t.type, *
[dim.value for dim in decl.init_dim])
6          symbol = VarSymbol(decl.ident.value, decl_type)
7          self.arrays.append(symbol)
8      else:
9          decl_type = decl.var_t.type
10         symbol = VarSymbol(decl.ident.value, decl_type)
11     ...
12
13 def visitIndexExpr(self, expr: IndexExpr, ctx: ScopeStack) ->
    None:
14     if isinstance(expr.base, Identifier) and not
ctx.lookupOverStack(expr.base.value):
15         raise DecafUndefinedVarError(expr.base.value)
16     expr.base.accept(self, ctx)
17     expr.index.accept(self, ctx)
18     #! 根据 base 类型设置 expr 的类型
19     if isinstance(expr.base, Identifier):
20         expr.setattr('type',
expr.base.getattr('symbol').type.indexed)
21     else:
22         expr.setattr('type', expr.base.getattr('type').indexed)
```

通过对节点的 `setattr('type')` 与 `getattr('type')` 操作实现类型运算一致性检查, 一共有 `INT` 与 `ArrayType` 两种类型, 如果相应运算类型不一致, 则会抛出异常 `DecafBadReturnTypeError`.

中间代码生成: 在 `utils/tac/tacgen.py` 中 `TACFuncEmitter` 类实现了依地址的读写数组元素的方法:

```
1 # frontend/tacgen/tacgen.py
2 class TACFuncEmitter(TACVisitor):
3     def visitLoadByAddress(self, addr: Temp):
4         dst = self.freshTemp()
5         self.func.add(LoadIntLiteral(dst, addr, 0))
6         return dst
7
8     def visitStoreByAddress(self, value: Temp, addr: Temp):
9         self.func.add(StoreIntLiteral(value, addr, 0))
```

修改 `TACGen` 类的 `visitIndexExpr`, `visitAssignment` 及 `visitIdentifier` 方法, 设置数组索引表达式的地址, 并判断访问标识符和赋值操作是否对数组进行:

```
1 # utils/tac/tacgen.py
2 class TACGen(Visitor[TACFuncEmitter, None]):
3     def visitIdentifier(self, ident: Identifier, mv:
4 TACFuncEmitter) -> None:
5         # 数组类型 -> 设置数组地址
6         if isinstance(ident.getattr('symbol').type, ArrayType):
7             ident.setattr('addr',
8 mv.visitLoadAddress(ident.getattr('symbol')))
9         ...
10
11     def visitIndexExpr(self, expr: IndexExpr, mv: TACFuncEmitter)
12 -> None:
13         expr.base.setattr('slice', True)
14         expr.base.accept(self, mv)
15         expr.index.accept(self, mv)
```

```

13         #! 递归计算当前索引的偏移量
14         addr = mv.visitLoad(expr.getattr('type').size)
15         mv.visitBinarySelf(tacop.TacBinaryOp.MUL, addr,
expr.index.getattr('val'))
16         mv.visitBinarySelf(tacop.TacBinaryOp.ADD, addr,
expr.base.getattr('addr'))
17         expr.setattr('addr', addr)
18         #! 递归完毕, 通过地址获得数组元素值
19         #! `slice` 属性表示为数组切片, 无需获取数据
20         #! 保证递归结束只有完整的索引表达式设置了返回值
21         if not expr.getattr('slice'):
22             expr.setattr('val', mv.visitLoadByAddress(addr))
23
24     def visitAssignment(self, expr: Assignment, mv:
TACFuncEmitter) -> None:
25         # 索引类型 -> 访问数组地址
26         if isinstance(expr.lhs, IndexExpr):
27             expr.lhs.setattr('slice', True)
28             expr.lhs.accept(self, mv)
29             mv.visitStoreByAddress(expr.rhs.getattr('val'),
expr.lhs.getattr('addr'))
30             expr.setattr('val', expr.rhs.getattr("val"))
31         ...

```

目标代码生成: 在 `utils/riscv.py` 中实现了 `addi` 的 RISC-V 指令, 用于加载通过 `TACFunc.arrays` 传递并保存在栈的局部数组 (仍通过 `Program.globalVars` 加载全局数组地址):

```

1  # utils/riscv.py
2  class ImmAdd(TACInstr):
3      def __init__(self, dst: Temp, src: Temp, value: int):
4          super().__init__(InstrKind.SEQ, [dst], [src], None)
5          self.value = value
6
7      def __str__(self) -> str:
8          assert -2048 <= self.value <= 2047 # Riscv imm [11:0]
9          return "addi " + Riscv.FMT3.format(
10              str(self.dsts[0]), str(self.srcs[0]), str(self.value)
11          )

```

进入一个函数前, 提前将其中的局部数组压栈, 在 `backend/subroutineinfo.py` 中计算得到各数组偏移量及占用栈帧大小:

```

1  # backend/subroutineinfo.py
2  class SubroutineInfo:
3      def __init__(self, funcLabel: FuncLabel, numArgs: int, arrays:
Dict[str, VarSymbol]) -> None:
4          self.offsets: Dict[str, int] = {}
5          self.size = 0
6
7          for name, symbol in self.arrays.items():
8              self.offsets[name] = self.size
9              self.size += symbol.type.size

```

在 `backend/riscv/riscvasmemitter.py` 中, `RiscvAsmEmitter` 对象初始化时打印全局数组至 `.bss` 区, `RiscvInstrSelector` 中通过偏移量实现局部数组访问:

```

1  # backend/riscv/riscvasmemitter.py
2  class RiscvAsmEmitter(AsmEmitter):
3      class RiscvInstrSelector(TACVisitor):
4          def visitLoadAddress(self, instr: LoadAddress) -> None:
5              if instr.symbol.isGlobal:
6                  ...
7              else:
8                  self.seq.append(Riscv.ImmAdd(instr.dsts[0],
Riscv.SP, self.info.offsets[instr.symbol.name]))

```

此时的 `self.nextLocalOffset` 及保存 `RA`, `FP` 和 `Callee-saved` 寄存器时需要额外加上 `self.info.size` 由局部数组占用的栈帧大小.

1.3 step 12

代码流读入: 在 `main.py` 中给出 `memset` 函数 `fill_csx` 的字符串表示, 命名方式是为了防止可能的函数重名 (虽然根据给出的测例, 没有名为 `fill_csx` 的函数, 但更好的方法是在语法分析后动态命名). 随后直接将代码加入输入的 `code` 前:

```

1  # main.py
2  memsetFunc = r"""
3  int fill_csx(int array[], int cnt) {
4      for (int i = 0; i < cnt; i = i + 1) {
5          array[i] = 0;
6      }
7      return 0;
8  }
9  """
10
11 def step_parse(args: argparse.Namespace):
12     code = memsetFunc + "\n" + readCode(args.input)
13     r: Program = parser.parse(code, lexer=lexer)
14     ...

```

词法语法分析: 修改 `frontend/ast/tree.py`, 修改 `Parameter` 节点定义, 添加维度声明 `init_dim` 标识数组传参; 同时添加数组初始化列表节点 `InitList`.


```

1  # frontend/ast/tree.py
2  class InitList(Node):
3      def __init__(self, init_list: List[IntLiteral]):
4          super().__init__("init_list")
5          self.init_list = init_list
6          self.value = [item.value for item in init_list]

```

在 `frontend/parser/ply_parser.py` 中参考已有实现, 给出相关 `CFG` 文法.

语义分析: 修改 `frontend/typecheck/namer.py`, 在 `visitFunction` 中记录函数声明中进行传参的参数数组; 在 `visitParameter` 中依据参数类型进行标识符 `symbol` 生成, 如果是数组则调用 `ArrayType.multidim` 生成相应的数据类型.

```

1  # frontend/typecheck/namer.py
2  def visitParameter(self, param: Parameter, ctx: ScopeStack) ->
    None:
3      ...
4      #! 为数组参数, 检查并生成相应的数据类型
5      if param.init_dim:
6          for index, dim in enumerate(param.init_dim):
7              if index == 0 and dim is NULL:
8                  continue
9              if dim.value <= 0:
10                 raise DecafBadArraySizeError()
11             decl_type = ArrayType.multidim(param.var_t.type, *
[dim.value if dim else None for dim in param.init_dim])
12             symbol = VarSymbol(param.ident.value, decl_type)
13             ...

```

中间代码生成: 在 `utils/tac/tacgen.py` 中修改 `TACGen` 类的 `visitIdentifier` 方法, 对全局数组与局部数组直接加载数组地址, 对参数数组加载相应虚拟寄存器; 修改 `visitDeclaration` 方法, 若为带初始化列表的局部数组声明, 先调用 `fill_csx` 函数进行内存清零, 然后逐一进行元素初始化:

```

1  # utils/tac/tacgen.py

```

```

2  class TACGen(Visitor[TACFuncEmitter, None]):
3      def visitIdentifier(self, ident: Identifier, mv:
TACFuncEmitter) -> None:
4          if isinstance(ident.getattr('symbol').type, ArrayType):
5              # 全局数组与局部数组
6              if ident.getattr("symbol").isGlobal or
ident.getattr('symbol') not in mv.func.p_arrays.values():
7                  ident.setattr('addr',
mv.visitLoadAddress(ident.getattr('symbol')))
8                  ident.setattr('val', ident.getattr('addr'))
9                  # 参数数组
10             else:
11                 ident.setattr('addr',
ident.getattr('symbol').temp)
12                 ident.setattr('val', ident.getattr('addr'))
13
14         def visitDeclaration(self, decl: Declaration, mv:
TACFuncEmitter) -> None:
15             if isinstance(decl.init_expr, InitList):
16                 #! 调用 `fill_csx` 函数进行初始化
17                 symbol = decl.getattr("symbol")
18                 addr = mv.visitLoadAddress(symbol)
19                 #! size 为 4 -> int 字长
20                 size = symbol.type.full_indexed.size
21                 interval = mv.visitLoad(size)
22                 mv.visitParam(addr)
23                 mv.visitParam(mv.visitLoad(symbol.type.size // size))
24                 mv.visitCall(FuncLabel("fill_csx"))
25                 #! 依次将初始化列表中的值存入数组中
26                 for value in decl.init_expr.value:
27                     mv.visitStoreByAddress(mv.visitLoad(value), addr)
28                     mv.visitBinarySelf(tacop.TacBinaryOp.ADD, addr,
interval)

```

目标代码生成: 在 `backend/riscv/riscvasmemitter.py` 中, `RiscvAsmEmitter` 对象初始化时打印带初始化列表的全局数组至 `.data` 区.

由于传参和调用分离, 对于参数数组的使用与普通函数参数并无区别, 这一步沿用 step 9 的实现即可, 后端无需进行其他工作.

2 思考题

2.1 step 10

1. 写出 `la v0, a` 这一 RiscV 伪指令可能会被转换成的指令组合 (两种即可).

答: 查阅 [The RISC-V Instruction Set Manual](#).

- non-PIC 可能转换为:

```
1 | auipc v0, delta[31:12] + delta[11]
2 | addi v0, v0, delta[11:0]
```

其中 `delta` 为 `a` 相对 `PC` 的偏移量, 因为地址在编译时已知, 所以使用 `addi` 加载.

- PIC 可能转换为:

```
1 | auipc v0, delta[31:12] + delta[11]
2 | lw v0, v0, delta[11:0]
```

其中 `delta` 为 `a` 在 GOT 中的地址相对 `PC` 的偏移量, 因为 GOT 中的地址可能在运行时进行重定位, 因此需要使用 `lw` 从内存中加载地址.

- 两种指令组合都存在 `+ delta[11]`, 这保证了在 `delta[11] = 1` 时, 低 12 位经过符号扩展, 最终也能够得到正确的结果.

2.2 step 11

1. C 语言规范规定, 允许局部变量是可变长度的数组 (VLA), 在我们的实验中为了简化, 选择不支持它. 请简要回答, 如果支持一维的可变长度的数组 (类似 `int n = 5; int a[n];`, 但不允许类似 `int n = ...; int m = ...; int a[n][m];`), 而且要求数组仍保存在栈上 (不允许用堆上动态内存申请), 应该在现有的实现基础上做出那些改动?

提示: 不能再在进入函数时统一给局部变量分配内存, 离开时统一释放内存.

答: 不能在进入函数时为 VLA 分配内存. 由于 VLA 的大小在编译期确定, 运行到声明 VLA 时, 将当前 `SP` 和 VLA 大小 `size` 保存到栈上 `-size(SP)` 处, 并移动 `SP`, 此时先前存储好的 `SP` 和 `size` 恰位于栈顶; 访问 VLA 的元素时, 计算偏移量 VLA 元素地址; 当 VLA 离开作用域时, 恢复 `SP` 即可.

2.3 step 12

1. 作为函数参数的数组类型第一维可以为空. 事实上, 在 C/C++ 中即使标明了第一维的大小, 类型检查依然会当作第一维是空的情况处理. 如何理解这一设计?

答: C/C++ 数组传参时, 数组名对应首元素地址; 函数无需为数组分配内存, 只需要通过首地址和偏移量即可访问到任意数组元素, 根据数组索引计算偏移量不会用到第一维的大小, 因此编译器会将数组参数的第一维视为空.