

LLVM 影响下的编译器新生态

高一 4 班 - 梁韬 - GitHub@Origami404

2019 年 2 月 14 日

1 前言

21 世纪以来, 新的编程语言的产生速度越来越快, 包括 DSL 在内的全新编程语言对于编译技术的质和量的要求也越来越高. 而反观 Intel 等硬件厂商, 指令集也越发复杂, 同时编译优化技术特别是 GC 及 JIT 技术的发展对于编译器作者的水平也提出了极高的要求. 面对这种困难, LLVM 诞生了.

LLVM 是一个**自由软件项目**, 是一套编译器基础设施, 基于 C++. LLVM 的命名最早源于所谓**底层虚拟机 (Low Level Virtual Machine)** 但随着其本身的发展及壮大, 现在的 LLVM 只是单纯的作为一个无意义的名称, 用于统称 LLVM 项目下产生的各种工具, 如 LLVM IR, lldg, LLVM C++ Standard Library 等.

本综述主要叙述 LLVM 对于编译器编写这一传统领域所带来的影响及基于 LLVM 编写编译器前端的基本流程.

2 LLVM 对传统编译编译器编写的影响

传统的编译器实现一般包括两大部分七大模块. 第一个部分是**分析**. 分析部分将源代码分解为多个组成要素, 并在这些组成要素上添加语法结构. 在此部分编译器将捕获并处理全部的词法与语法错误, 收集整理与代码有关的其他信息, 最后构造出一种**中间表示 (IR)**.

第二个部分是**综合**. 在此部分编译器开始着手优化并产生目标代码. 在这一部分编译器可能使用某些巧妙的算法进行机器无关的代码优化, 又或者使用某种依赖于目标平台或目标语言特性的优化机理进行优化, 最后产生的便是目标代码.

在此划分之上编译器又被按功能分为七个模块, 分别是**词法分析**, **语法分析**, **语义分析**, **中间代码生成**, **机器无关代码优化**, **代码生成**, **机器相关代码生成**. 其中词法分析, 语法分析, 语义分析及中间代码生成属于分析部分; 机器无关代码优化, 代码生成, 机器相关代码生成属于综合部分.

读者不难发现, IR 作为联系前端和后端的形式, 在编译器编写中具有重要作用. 在一门语言的实现中, 除了语言本身的语法定义, 就属 IR 的定义最为重要. 泛观世界上成熟的编译器工具集和 VM 如 GCC, LLVM, JVM, 甚至是 LuaVM 都定义了极为成熟和在一定程度上泛用的 IR 表示. 这些 IR 为平台的二次开发奠定了基础, 也是其能流行于世界的重要原因. 而相比起其他编译平台, LLVM 的优势在于其 IR 的开发一开始就基于一个目标——多语言通用, 多范式支持, 多用途使用. 因而在快速实现编译器时 LLVM 具有天然的优势, 而其完全透明, 类型完备的 IR 也为基于该语言的开发工具提供了实现包括语法高亮, 语义代码补全等面给用户的功能的极大便利. 诸如 Apple 的 Xcode, Microsoft 的 VSCode 的 C++ 插件等都支持对 LLVM 编译出的 IR 进行利用而改善用户体验的功能. 这也是 LLVM 的一大优势.

3 利用 LLVM 实现编译器前端的基本流程

LLVM C++ 库对中间代码生成, 优化及跨平台机器代码生成提供了简洁易用的接口. 在一般的编写流程中, 开发人员主要编写词法与语法分析器, 即俗称的 Scanner 与 Parser, 而后通过在 Parser 得出的 AST 继承体系中加入调用 LLVM 生成 IR 的方法来生成 LLVM IR, 极大的简化了开发流程. 对于某些对速度要求不高的需求, 开发人员甚至可以使用诸如 Bison 和 flex 等分析器生成器生成语法词法分析器, 从而极大地降低了编译器开发的门槛. 同时也因为 LLVM IR 的生成一般基于 AST, 而在 AST 层面往往可以对语言作出极其灵活的修改而不改变用户代码的感知, 故这种方法也具有极高的灵活性, 甚至在某些情况下能实现超出 Python 的 metaclass 和 Java 的 Reflection 的令人惊奇的灵活修改, 能在源代码层面实现原本许多嵌入式语言中在宿主层面才实现的特性.

为更好的介绍 LLVM C++ 接口, 在这里实现一个简短的带变量二元运算解析. 我们的目标是支持诸如:

```
1  def x = 2.44
2  def y_s = 300.45 + x
3  def c = (x + y_s)*3
4  c * 2 + x / y_s
```

的变量赋值与二元运算. 二元运算支持 +, -, *, / 与 (), 为了简便所有量的类型均为 double. 更严谨的, 此语言的 BNF 如下: (也是一个 LL(1) 文法)

$$programmer \rightarrow define \mid expression \quad (1)$$

$$define \rightarrow \mathbf{def} \text{ Id } = expression \quad (2)$$

$$expression \rightarrow primary [binop] \quad (3)$$

$$primary \rightarrow Id \mid Num \mid (expression) \quad (4)$$

$$binop \rightarrow primary \text{ Op } binop \quad (5)$$

由于不是重点, 这里直接给出 AST 以及 Scanner 和 parser 的定义:

```
1 // 编译命令
2 // clang++ -g -O3 *.cpp 'llvm-config --cxxflags --ldflags --system-libs --libs core'
3 #include "llvm/ADT/APFloat.h"
4 #include "llvm/ADT/STLExtras.h"
5 #include "llvm/IR/Constants.h"
6 #include "llvm/IR/DerivedTypes.h"
7 #include "llvm/IR/IRBuilder.h"
8 #include "llvm/IR/LLVMContext.h"
9 #include "llvm/IR/Module.h"
10 #include "llvm/IR/Type.h"
11 #include "llvm/IR/Verifier.h"
12 #include <algorithm>
13 #include <cctype>
14 #include <cstdio>
15 #include <cstdlib>
16 #include <map>
17 #include <memory>
18 #include <string>
19 #include <vector>
```

```

20 using std::unique_ptr;
21 using std::string;
22 using std::map;
23 using std::move;
24 using namespace llvm;

```

头文件

```

1 // 这里是AST节点的定义
2 namespace {
3     class ExprAST {
4     public:
5         virtual ~ExprAST() {}
6         // 用于生成IR
7         virtual Value* codegen() = 0;
8     };
9     using ASTptr = unique_ptr<ExprAST>;
10    class NumAST : public ExprAST {
11    public:
12        double val;
13        NumAST(double v): val(v) {}
14        virtual Value* codegen();
15    };
16    class VarAST : public ExprAST {
17    public:
18        string name;
19        VarAST(const string& n): name(n) {}
20        virtual Value* codegen();
21    };
22    class BinOpAST : public ExprAST {
23    public:
24        char op;
25        unique_ptr<ExprAST> lhs, rhs;
26        BinOpAST(char op, ASTptr LHS, ASTptr RHS)
27            : op(op), lhs(move(LHS)), rhs(move(RHS)) {}
28        virtual Value* codegen();
29    };
30    class DefAST {
31    public:
32        string name;
33        ASTptr val;
34        DefAST(const string& name, ASTptr val)
35            : name(name), val(move(val)) {}
36        virtual void codegen();
37    };
38 }

```

AST

下面是一些主要环境变量和辅助函数:

```

1 enum Token {
2     tok_eof = -1,
3     tok_identifier = -2,
4     tok_number = -3,
5     tok_def = -4,
6     tok_end = -5
7 };

```

```

8
9 static string identifier_str;
10 static double num_val;
11
12 int gettok() {
13     static int last_char = ' ';
14     while (isspace(last_char))
15         last_char = getchar();
16     // identifier || def
17     if (isalpha(last_char)) {
18         identifier_str = last_char;
19         while (isalnum(last_char = getchar()))
20             identifier_str += last_char;
21         // def
22         if (identifier_str == "def")
23             return tok_def;
24         return tok_identifier;
25     }
26     // Number, a small bug
27     if (isdigit(last_char) || last_char == '.') {
28         string num_str;
29         do {
30             num_str += last_char;
31             last_char = getchar();
32         } while (isdigit(last_char) || last_char == '.');
33         num_val = strtod(num_str.c_str(), nullptr);
34         return tok_number;
35     }
36     // EOF
37     if (last_char == EOF)
38         return tok_eof;
39     if (last_char == ';') {
40         last_char = getchar();
41         return tok_end;
42     }
43
44     // Operators
45     int this_char = last_char;
46     last_char = getchar();
47     return this_char;
48 }
49 // 递归下降解析器的预读字符缓冲
50 int cur_tok;
51 int next_token() {
52     return (cur_tok = next_token());
53 }
54 ASTptr Error(const char* str) {
55     fprintf(stderr, "Error: %s\n", str);
56     return nullptr;
57 }
58 // 用于递归下降解析的运算符优先级表
59 map<char, int> op_prec;
60 int get_prec(char c) {
61     auto prec = op_prec.find(c);
62     if (prec == op_prec.end())
63         return -1;
64     else return prec->second;

```

```

65 }
66
67 ASTptr Error(const char* str) {
68     fprintf(stderr, "Error: %s\n", str);
69     return nullptr;
70 }

```

环境

同时一个简单的解析器也在下方给出, 注意本代码出于演示目的忽略了很多错误处理.

```

1  ASTptr parse_expression();
2  unique_ptr<DefAST> parse_define();
3  ASTptr parse_primary();
4  ASTptr parse_binop(int expr_prec, ASTptr lhs);
5  ASTptr parse_id();
6  ASTptr parse_num();
7  ASTptr parse_paren();
8
9  /// programmer ::= define | expression
10 // 同时承担程序入口的责任
11 void parse_programmer() {
12     while (true) {
13         printf("ready> ");
14         switch (cur_tok) {
15             case tok_eof: return;
16             case tok_def: parse_define()->codegen(); break;
17             case tok_end: break;
18             default: parse_expression()->codegen()->print(errs());
19         }
20         next_token();
21     }
22 }
23 /// define ::= 'def' Id '=' expression
24 unique_ptr<DefAST> parse_define() {
25     next_token(); // eat def
26     auto name = identifier_str;
27     next_token(); // eat '='
28     next_token();
29     auto val = parse_expression();
30     return llvm::make_unique<DefAST>(DefAST(name, move(val)));
31 }
32 /// expr ::= primary [binop]
33 ASTptr parse_expression() {
34     return parse_binop(0, parse_primary());
35 }
36 /// primary ::= Id | Num | paren
37 ASTptr parse_primary() {
38     switch (cur_tok) {
39         case tok_identifier:
40             return parse_id();
41         case tok_number:
42             return parse_num();
43         case '(':
44             return parse_paren();
45         default:
46             return Error("Unknown token in primary");
47     }

```

```

48 }
49 /// paren ::= '(' expression ')'
50 ASTPtr parse_paren() {
51     next_token(); // eat (
52     auto expr = parse_expression();
53     if (cur_tok != ')')
54         return Error("Unmatch paren");
55     next_token(); // eat )
56     return expr;
57 }
58 /// binop ::= primary Op binop
59 ASTPtr parse_binop(int expr_prec, ASTPtr lhs) {
60     while (true) {
61         int now_prec = get_prec(cur_tok);
62         if (now_prec < expr_prec)
63             return lhs;
64
65         int binop = cur_tok;
66         next_token();
67         auto rhs = parse_primary();
68
69         int next_prec = get_prec(cur_tok);
70         if (now_prec < next_prec)
71             rhs = parse_binop(expr_prec + 1, move(rhs));
72         lhs = llvm::make_unique<BinOpAST>(binop, move(lhs), move(rhs));
73     }
74 }
75
76 ASTPtr parse_num() {
77     auto res = llvm::make_unique<NumAST>(num_val);
78     next_token();
79     return move(res);
80 }
81
82 ASTPtr parse_id() {
83     auto res = llvm::make_unique<VarAST>(identifier_str);
84     next_token();
85     return move(res);
86 }

```

解析器

下面就是本节讲述的重点: LLVM IR C++ API. 我们使用 AST 节点的 `codegen` 虚方法来递归地生成 IR, 您将看到这种生成方法是如此强大与方便, 以至于代码量还不及解析器:

```

1 // LLVM API 文档(英文): https://llvm.org/docs/LangRef.html
2 static LLVMContext my_context;
3 static IRBuilder<> builder(my_context);
4 static unique_ptr<Module> my_module;
5 static map<string, Value*> named_vals;
6 Value *LogErrorV(const char *str) {
7     LogError(str);
8     return nullptr;
9 }
10
11 Value* NumAST::codegen() {
12     return ConstantFP::get(my_context, APFloat(this->val));
13 }

```

```

14 Value* VarAST::codegen() {
15     Value* v = named_vals[this->name];
16     if (!v)
17         return LogErrorV("undefine variable name");
18     return v;
19 }
20 Value* BinOpAST::codegen() {
21     Value* l = this->lhs->codegen();
22     Value* r = this->rhs->codegen();
23     if (!l || !r)
24         return nullptr;
25     switch (this->op) {
26         case '+':
27             return builder.CreateFAdd(l, r, "addtmp");
28         case '-':
29             return builder.CreateFSub(l, r, "subtmp");
30         case '*':
31             return builder.CreateFMul(l, r, "multmp");
32         case '/':
33             return builder.CreateFDiv(l, r, "divtmp");
34         default: return LogErrorV("unexpected binop");
35     }
36 }
37 void DefAST::codegen() {
38     printf("Name: %s defined\n", this->name.c_str());
39     named_vals[this->name] = this->val->codegen();
40 }

```

IR 生成

从上述代码可以看出 LLVM IR C++ API 的完善与简洁. 在保留了比 C 更贴近一点底层的表示的同时很好地维护了 IR 的类型信息. 完整的 LLVM API 可以在 LLVM 的网站上查询.

使用测试样例运行程序:

```

1 ready> def x = 42.0;
2 ready> def y = 3.2*x;
3 ready> x+y;
4 double 1.764000e+02
5 ready>

```

测试

可以看到 LLVM 会在生成 IR 中就为我们完成简单的常量常量折叠, 事实上在这个小例子中 LLVM 甚至能帮我们完成全部的运算并折叠成一个常量. 注意这在 LLVM 的分类中 “连优化都算不上”, 是默认开启的. LLVM 对机器无关的代码优化能力可见一斑.

但即使是这种小例子也足以折射出 LLVM IR 的一个特点: 强类型. LLVM IR 对于类型是极为严格的, 甚至不允许 C 的很多符合常理的隐式类型转换, 这一点也在后续的优化中发挥很大的作用.

4 总结与参考引用

本综述阐述了基于 LLVM 编写编译器前端的基本流程, 为今后开展编译技术方面研究性学习奠定了坚实的基础. 今后的研究性学习将着手于编译器前端技术的学习, 实现一个简单的 Scheme 解释器与 C 编译器.

本文主要参考: LLVM Tutorial : <https://llvm.org/docs/tutorial/index.html>

LLVM Document : <https://llvm.org/docs/LangRef.html>

Wikipedia : <https://www.wikipedia.org/>

« 编译原理 » 机械工业出版社, 第二版.

« 深入理解计算机系统 » 机械工业出版社, 第三版.

代码地址:

<https://github.com/Origami404/Study/blob/master/cal.cpp>