网络空间安全学院

编译系统原理实验报告

# 了解编译器及 LLVM IR 编程

姓名：王天鹏

学号：2013013

专业：信息安全与法学双学位班

2022 年 9 月

# 目录

# 1 实验背景介绍

## 1.1 问题描述

　　本次实验的基本研究对象为编译器的基本构成和工作原理。笔者以一个简单的 C 语言竭诚的源程序为例，通过调整编译器的程序选项获得各阶段的输出，研究它们与源程序的关系，从而进一步了解一个完整的编译过程，熟悉预处理器、编译器、汇编器、链接器在将程序翻译成汇编语言再到二进制文件过程中各自的分工。本次实验所使用的阶乘程序的 C 语言源代码如下：

```c
#include <stdio.h>
#define N 200
    #ifndef N
        scanf("%d", &n);
#else
#pragma message ("N is defined\n")
#endif
int main()
{
    int i, n, f;
    scanf("%d", &i);
    i = 2;
    f = 1;
    while (i <= n)
    {
        f = f * i;
        i = i + 1;
    }
    printf("f= %d\n", f);
    return 0;
}

```

## 1.2 实验环境

　　在本次实验中，笔者使用了 Linux 云服务器，并通过 xshell 进行远程连接，使用的软件版本如下

1. 操作系统：Linux VM-4-12-centos 3.10.0-1160.66.1.el7.x86_64 #1 SMP x86_64 GNU/Linux

2. GCC 编译器：gcc (GCC) 4.8.5 20150623 (Red Hat 4.8.5-44)

3. CLANG 前端：clang version 3.4.2 (tags/RELEASE_34/dot2-final)

4. LLVM 后端：14.0.0

5. GCC-ARM 编译器：arm-linux-gnueabihf-gcc 11.2.0

6. QEMU-ARM 模拟器：6.2.0

## 2   预处理器做了什么？

　　当源程序被编译器处理后，第一个环节是预处理器先处理。预处理器的工作主要是处理源代码中以 开始的预编译指令，例如展开所有宏定义、插入 include 指向的文件等，以获得经过预处理的源程序。通过命令 g 始的预编译指令，以及所有宏定义进行插入 include 指向的文件和宏定义语句。其中预处理后文件 main.i 文件有几千行之多，下列仅对 main.i 文件进行了部分展示，省略了大篇幅的对于 stdio.h 头文件的引用部分的展示。

```
1
2
3    # 1 "main.c"
4    1 "<built-in>"
5    1 "<command-line>"
6    1 "/usr/include/stdc-predef.h" 1 3 4
7    1 "<command-line>" 2
8    1 "main.c"
9    1 "/usr/include/stdio.h" 1 3 4
10   27 "/usr/include/stdio.h" 3 4
11   1 "/usr/include/features.h" 1 3 4
12   375 "/usr/include/features.h" 3 4
13   1 "/usr/include/sys/cdefs.h" 1 3 4
14   392 "/usr/include/sys/cdefs.h" 3 4
15   1 "/usr/include/bits/wordsize.h" 1 3 4
16   393 "/usr/include/sys/cdefs.h" 2 3 4
17   376 "/usr/include/features.h" 2 3 4
18   399 "/usr/include/features.h" 3 4
19   1 "/usr/include/gnu/stubs.h" 1 3 4
20   10 "/usr/include/gnu/stubs.h" 3 4
21   1 "/usr/include/gnu/stubs-64.h" 1 3 4
22   11 "/usr/include/gnu/stubs.h" 2 3 4
23   400 "/usr/include/features.h" 2 3 4
24   28 "/usr/include/stdio.h" 2 3 4
25   此处省略了对于 stdio.h 头文件的替换部分。
26
27   int main()
28   {
29       int i, n, f;
30       scanf("%d", &i);
31       i = 2;
32       f = 1;
33       while (i <= n)
34       {
```

```
35          f = f * i;
36          i = i + 1;
37      }
38      printf("f= %d\n", f);
39      return 0;
40  }
41
```

如上述代码，可以看到，

- 预处理器对源程序中的宏定义N了替换，将N替换为了常量 200。

- 源程序中的注释在这一阶段被删除了。

- 条件编译语句如#idndef、#else被预处理器进行删除替换。

- 预处理器添加了行号和文件名等内容。

- 值得注意的是，gcc 的预处理器并未对于#pragma once进行处理。这是由于在 C/C++ 中，#pragma once是一个非标准但是被广泛支持的方式。#pragma once方式产生于#ifndef之后。#ifndef方式受 C/C++ 语言标准的支持，不受编译器的任何限制；而#pragma once方式有些编译器不支持，兼容性不够好，本次实验的 gcc 编译器预处理器对于#pragma once未作处理。#ifndef可以针对一个文件中的部分代码，而#pragma once只能针对整个文件。相对而言，#ifndef更加灵活，兼容性好，#pragma once操作简单，效率高。

# 3  编译器做了什么？

## 3.1  词法分析

经过预处理器处理后，编译器会对预处理后的 main.i 进行词法分析。通过lang -E -Xlang -dump-tokens main.c 选项，编译器会将源程序转换为单词序列，因此我们可以获得文本格式的 token 序列，阶乘程序的部分 tokens 序列如下所示：

```
1   typedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/bin/../lib/clang/3.4.2/include/stddef.h:34:1>
2  ng 'long'^^I [LeadingSpace]^^ILoc=</usr/bin/../lib/clang/3.4.2/include/stddef.h:34:9 <Spelling=<
3  t 'int'^^I [LeadingSpace]^^ILoc=</usr/bin/../lib/clang/3.4.2/include/stddef.h:34:9 <Spelling=<bu
4  entifier 'ptrdiff_t'^^I [LeadingSpace]^^ILoc=</usr/bin/../lib/clang/3.4.2/include/stddef.h:34:26
5  mi ';'^^I^^ILoc=</usr/bin/../lib/clang/3.4.2/include/stddef.h:34:35>
6  pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/bin/../lib/clang/3.4.2/include/stddef.h:42:1>
7  ng 'long'^^I [LeadingSpace]^^ILoc=</usr/bin/../lib/clang/3.4.2/include/stddef.h:42:9 <Spelling=<
8  signed 'unsigned'^^I [LeadingSpace]^^ILoc=</usr/bin/../lib/clang/3.4.2/include/stddef.h:42:9 <Sp
9  t 'int'^^I [LeadingSpace]^^ILoc=</usr/bin/../lib/clang/3.4.2/include/stddef.h:42:9 <Spelling=<bu
10 entifier 'size_t'^^I [LeadingSpace]^^ILoc=</usr/bin/../lib/clang/3.4.2/include/stddef.h:42:23>
11 mi ';'^^I^^ILoc=</usr/bin/../lib/clang/3.4.2/include/stddef.h:42:29>
12 pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/bin/../lib/clang/3.4.2/include/stddef.h:65:1>
```

13 t 'int'^^I [LeadingSpace]^^ILoc=</usr/bin/../lib/clang/3.4.2/include/stddef.h:65:9 <Spelling=<bu

14 entifier 'wchar_t'^^I [LeadingSpace]^^ILoc=</usr/bin/../lib/clang/3.4.2/include/stddef.h:65:24>

15 mi ';'^^I^^ILoc=</usr/bin/../lib/clang/3.4.2/include/stddef.h:65:31>

16 pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/include/bits/types.h:30:1>

17 signed 'unsigned'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:30:9>

18 ar 'char'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:30:18>

19 entifier '__u_char'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:30:23>

20 mi ';'^^I^^ILoc=</usr/include/bits/types.h:30:31>

21 pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/include/bits/types.h:31:1>

22 signed 'unsigned'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:31:9>

23 ort 'short'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:31:18>

24 t 'int'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:31:24>

25 entifier '__u_short'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:31:28>

26 mi ';'^^I^^ILoc=</usr/include/bits/types.h:31:37>

27 pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/include/bits/types.h:32:1>

28 signed 'unsigned'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:32:9>

29 t 'int'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:32:18>

30 entifier '__u_int'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:32:22>

31 mi ';'^^I^^ILoc=</usr/include/bits/types.h:32:29>

32 pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/include/bits/types.h:33:1>

33 signed 'unsigned'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:33:9>

34 ng 'long'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:33:18>

35 t 'int'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:33:23>

36 entifier '__u_long'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:33:27>

37 mi ';'^^I^^ILoc=</usr/include/bits/types.h:33:35>

38 pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/include/bits/types.h:36:1>

39 gned 'signed'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:36:9>

40 ar 'char'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:36:16>

41 entifier '__int8_t'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:36:21>

42 mi ';'^^I^^ILoc=</usr/include/bits/types.h:36:29>

43 pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/include/bits/types.h:37:1>

44 pelling=/usr/include/bits/types.h:103:35>>

45 t 'int'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:145:12 <Spelling=/usr/include/bits/ty

46 entifier '__rlim_t'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:145:26>

47 mi ';'^^I^^ILoc=</usr/include/bits/types.h:145:34>

48 pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/include/bits/types.h:146:1 <Spelling=/usr/include/b

49 signed 'unsigned'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:146:12 <Spelling=/usr/includ

50 ng 'long'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:146:12 <Spelling=/usr/include/bits/

51 t 'int'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:146:12 <Spelling=/usr/include/bits/ty

52 entifier '__rlim64_t'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:146:28>

53 mi ';'^^I^^ILoc=</usr/include/bits/types.h:146:38>

54 pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/include/bits/types.h:147:1 <Spelling=/usr/include/b

```
55 signed 'unsigned'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:147:12 <Spelling=/usr/inclu
56 t 'int'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:147:12 <Spelling=/usr/include/bits/ty
57 entifier '__id_t'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:147:24>
58 mi ';'^^I^^ILoc=</usr/include/bits/types.h:147:30>
59 pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/include/bits/types.h:148:1 <Spelling=/usr/include/b
60 ng 'long'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:148:12 <Spelling=/usr/include/bits/
61 t 'int'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:148:12 <Spelling=/usr/include/bits/ty
62 entifier '__time_t'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:148:26>
63 mi ';'^^I^^ILoc=</usr/include/bits/types.h:148:34>
64 pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/include/bits/types.h:149:1 <Spelling=/usr/include/b
65 signed 'unsigned'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:149:12
66 ng 'long'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:167:12 <Spelling=/usr/include/bits/
67 t 'int'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:167:12
68 t 'int'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:172:12 <Spelling=/usr/include/bits/ty
69 entifier '__fsblkcnt64_t'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:172:32>
70 mi ';'^^I^^ILoc=</usr/include/bits/types.h:172:46>
71 pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/include/bits/types.h:175:1 <Spelling=/usr/include/b
72 signed 'unsigned'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:175:12 <Spelling=/usr/inclu
73 ng 'long'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:175:12 <Spelling=/usr/include/bits/
74 t 'int'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:175:12 <Spelling=/usr/include/bits/ty
75 entifier '__fsfilcnt_t'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:175:30>
76 mi ';'^^I^^ILoc=</usr/include/bits/types.h:175:42>
77 pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/include/bits/types.h:176:1 <Spelling=/usr/include/b
78 signed 'unsigned'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:176:12 <Spelling=/usr/inclu
79 ng 'long'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:176:12 <Spelling=/usr/include/bits/
80 t 'int'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:176:12 <Spelling=/usr/include/bits/ty
81 entifier '__fsfilcnt64_t'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:176:32>
82 mi ';'^^I^^ILoc=</usr/include/bits/types.h:176:46>
83 pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/include/bits/types.h:179:1 <Spelling=/usr/include/b
84 ng 'long'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:179:12 <Spelling=/usr/include/bits/
85 t 'int'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:179:12 <Spelling=/usr/include/bits/ty
86 entifier '__fsword_t'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:179:28>
87 mi ';'^^I^^ILoc=</usr/include/bits/types.h:179:38>
88 pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/include/bits/types.h:181:1 <Spelling=/usr/include/b
89 ng 'long'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:181:12 <Spelling=/usr/include/bits/
90 t 'int'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:181:12 <Spelling=/usr/include/bits/ty
91 entifier '__ssize_t'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:181:27>
92 mi ';'^^I^^ILoc=</usr/include/bits/types.h:181:36>
93 pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/include/bits/types.h:184:1 <Spelling=/usr/include/b
94 ng 'long'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:184:12 <Spelling=/usr/include/bits/
95 t 'int'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:184:12 <Spelling=/usr/include/bits/ty
96 entifier '__syscall_slong_t'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:184:33>
```

```
97  mi ';'^^I^^ILoc=</usr/include/bits/types.h:184:50>
98  pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/include/b
99  signed 'unsigned'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:186:12 <Spelling=/usr/inclu
100 ng 'long'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:186:12 <Spelling=/usr/include/bits/t
101 t 'int'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:186:12
102 t 'int'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:195:12 <Spelling=/usr/include/bits/typ
103 entifier '__intptr_t'^^I [LeadingSpace]^^ILoc=</usr/include/bits/types.h:195:25>
104 mi ';'^^I^^ILoc=</usr/include/bits/types.h:195:35>
105 pedef 'typedef'^^I [StartOfLine]^^ILoc=</usr/include/bits/types.h:198:1 <Spelling=/usr/include/b
106 中间有省略
```

由于输出的 token 较长，文中仅展示了一小部分的 token 序列。通过上述 token，我们可以看出词法分析的任务是从左至右逐个字符地扫描源程序，产生一个个单词符号，把源程序的字符串改造成为单词符号串。clang 编译器的词法分析器会对词法单元、关键字、空格、缩进、换行、注释以及部分词法错误进行标记处理。通过对程序进行分词和标记，即可以将源代码转换成 tokens，以便于后续的语法处理。

## 3.2  语法分析

编译器会将词法分析生成的词法单元构建抽象语法树 (Abstract Syntax Tree，即 AST)。对于 gcc，笔者使用 gcc -fdump-tree-original-raw flag 获得文本格式的 AST 输出。LLVM 可以通过如下命令通过选项，我们获得了文本格式的 AST。

其中使用命令cc -fdump-tree-oriinal-raw main. 获得的 gcc 编译器的文本格式的 AST 输出如下，输出的结果保存在 main.c.003t.original 中。

```
1      ;; Function main (null)
2    enabled by -tree-original
3
4        bind_expr      type: @2      vars: @3      body: @4
5        void_type      name: @5      algn: 8
6        var_decl       name: @6      type: @7      scpe: @8
7                       srcp: main.c:4                size: @9
8                       algn: 32      used: 1
9        statement_list 0   : @10     1   : @11     2   : @12
10                       3   : @13     4   : @14     5   : @15
11                       6   : @16     7   : @17     8   : @18
12                       9   : @19     10  : @20     11  : @21
13                       12  : @22     13  : @23     14  : @24
14       type_decl      name: @25     type: @2
15       identifier_node strg: i      lngt: 1
16       integer_type   name: @26     size: @9      algn: 32
17                       prec: 32     sign: signed  min : @27
```

```
18                              max : @28
19        function_decl         name: @29         type: @30         srcp: main.c:2
20                              link: extern
21        integer_cst           type: @31         low : 32
22  0     decl_expr             type: @2
23  1     decl_expr             type: @2
24  2     decl_expr             type: @2
25  3     call_expr             type: @7          fn  : @32         0   : @33
26                              1   : @34
27  4     modify_expr           type: @7          op 0: @3          op 1: @35
28  5     modify_expr           type: @7          op 0: @36         op 1: @37
29  6     goto_expr             type: @2          labl: @38
30  7     label_expr            type: @2          name: @39
31  8     modify_expr           type: @7          op 0: @36         op 1: @40
32  由于篇幅原因，中间 18-81 行内容省略。
33  1     tree_list             valu: @58
34  2     type_decl             name: @88         type: @83
35  3     integer_type          name: @82         size: @51         algn: 8
36                              prec: 8           sign: signed      min : @84
37                              max : @85
38  4     integer_cst           type: @83         high: -1          low : -128
39  5     integer_cst           type: @83         low : 127
40  6     array_type            size: @89         algn: 8           elts: @83
41                              domn: @90
42  7     array_type            size: @91         algn: 8           elts: @83
43                              domn: @92
44  8     identifier_node       strg: char        lngt: 4
45  9     integer_cst           type: @31         low : 24
46  0     integer_type          size: @68         algn: 64          prec: 64
47                              sign: signed      min : @93         max : @94
48  1     integer_cst           type: @31         low : 56
49  2     integer_type          size: @68         algn: 64          prec: 64
50                              sign: signed      min : @93         max : @95
51  3     integer_cst           type: @96         low : 0
52  4     integer_cst           type: @96         low : 2
53  5     integer_cst           type: @96         low : 6
54  6     integer_type          name: @97         size: @68         algn: 64
55                              prec: 64          sign: unsigned min : @93
56                              max : @98
57  7     identifier_node       strg: sizetype lngt: 8
58  8     integer_cst           type: @96         low : -1
59
```

60

笔者同时使用命令`lang -E -Xlang -ast-dump main.c` 获得的 clang 编译器的文本格式的 AST 输出如下。

```
1  anslationUnitDecl 0x1c64d60 <<invalid sloc>>
2  TypedefDecl 0x1c65260 <<invalid sloc>> __int128_t '__int128'
3  TypedefDecl 0x1c652c0 <<invalid sloc>> __uint128_t 'unsigned __int128'
4  TypedefDecl 0x1c65610 <<invalid sloc>> __builtin_va_list '__va_list_tag [1]'
5  TypedefDecl 0x1c65670 </usr/bin/../lib/clang/3.4.2/include/stddef.h:34:1, col:26> ptrdiff_t 'long
6  TypedefDecl 0x1c656d0 <line:42:1, col:23> size_t 'unsigned long'
7  TypedefDecl 0x1c65730 <line:65:1, col:24> wchar_t 'int'
8  TypedefDecl 0x1c65790 </usr/include/bits/types.h:30:1, col:23> __u_char 'unsigned char'
9  TypedefDecl 0x1c657f0 <line:31:1, col:28> __u_short 'unsigned short'
10 TypedefDecl 0x1c65850 <line:32:1, col:22> __u_int 'unsigned int'
11 TypedefDecl 0x1c658b0 <line:33:1, col:27> __u_long 'unsigned long'
12 TypedefDecl 0x1c65910 <line:36:1, col:21> __int8_t 'signed char'
13 TypedefDecl 0x1c65970 <line:37:1, col:23> __uint8_t 'unsigned char'
14 TypedefDecl 0x1c659d0 <line:38:1, col:26> __int16_t 'short'
15 TypedefDecl 0x1c65a30 <line:39:1, col:28> __uint16_t 'unsigned short'
16 TypedefDecl 0x1cd5450 <line:40:1, col:20> __int32_t 'int'
17
18 `-NoThrowAttr 0x1d01cd0 </usr/include/sys/cdefs.h:56:35>
19 FunctionDecl 0x1d01ee0 </usr/include/stdio.h:356:12> fprintf 'int (FILE *, const char *, ...)' e
20 |-ParmVarDecl 0x1d01f80 <<invalid sloc>> 'FILE *'
21 |-ParmVarDecl 0x1d01fe0 <<invalid sloc>> 'const char *'
22 `-FormatAttr 0x1d02050 <col:12> printf 2 3
23 FunctionDecl 0x1d020a0 prev 0x1d01ee0 <col:1, line:357:43> fprintf 'int (FILE *, const char *, .
24 |-ParmVarDecl 0x1d01d20 <line:356:21, col:38> __stream 'FILE *restrict'
25 |-ParmVarDecl 0x1d01d90 <line:357:7, col:30> __format 'const char *restrict'
26 `-FormatAttr 0x1d02180 <line:356:12> printf 2 3
27 FunctionDecl 0x1d022c0 <line:362:12> printf 'int (const char *, ...)' extern
28 |-ParmVarDecl 0x1d02360 <<invalid sloc>> 'const char *'
29 `-FormatAttr 0x1d023c0 <col:12> printf 1 2
30 FunctionDecl 0x1d02410 prev 0x1d022c0 <col:1, col:56> printf 'int (const char *, ...)' extern
31 |-ParmVarDecl 0x1d021c0 <col:20, col:43> __format 'const char *restrict'
32 `-FormatAttr 0x1d024e0 <col:12> printf 1 2
33 FunctionDecl 0x1d02740 <line:364:12> sprintf 'int (char *, const char *, ...)' extern
34 |-ParmVarDecl 0x1d027e0 <<invalid sloc>> 'char *'
35 |-ParmVarDecl 0x1d02840 <<invalid sloc>> 'const char *'
36 `-FormatAttr 0x1d028b0 <col:12> printf 2 3
37 FunctionDecl 0x1d02900 prev 0x1d02740 <col:1, /usr/include/sys/cdefs.h:57:49> sprintf 'int (char
```

```
38 |-ParmVarDecl 0x1d02520 </usr/include/stdio.h:364:21, col:38> __s 'char *restrict'
39 |-ParmVarDecl 0x1d02590 <line:365:7, col:30> __format 'const char *restrict'
40 |-NoThrowAttr 0x1d029b0 </usr/include/sys/cdefs.h:57:37>
41 `-FormatAttr 0x1d029f0 </usr/include/stdio.h:364:12> printf 2 3
42 FunctionDecl 0x1d02c20 <line:371:12> vfprintf 'int ()' extern
43 `-FormatAttr 0x1d02cc0 <col:12> printf 2 0
44 FunctionDecl 0x1d02d10 prev 0x1d02c20 <col:1, line:372:24> vfprintf 'int (FILE *restrict, const c
45 |-ParmVarDecl 0x1d02a30 <line:371:22, col:39> __s 'FILE *restrict'
46 |-ParmVarDecl 0x1d02aa0 <col:44, col:67> __format 'const char *restrict'
47 |-ParmVarDecl 0x1d02b10 </usr/include/_G_config.h:46:20, /usr/include/stdio.h:372:19> __arg '__va
48 `-FormatAttr 0x1d02df0 <line:371:12> printf 2 0
49 FunctionDecl 0x1d02ff0 <line:377:12> vprintf 'int (const char *, __va_list_tag *)' extern
50
51
52 FunctionDecl 0x1d04440 prev 0x1d041b0 <col:1, line:392:62> vsnprintf 'int (char *, unsigned long
53 |-ParmVarDecl 0x1d03e80 <line:390:23, col:40> __s 'char *restrict'
54 |-ParmVarDecl 0x1d03ef0 <col:45, col:52> __maxlen 'size_t':'unsigned long'
55 |-ParmVarDecl 0x1d03f60 <line:391:9, col:32> __format 'const char *restrict'
56 |-ParmVarDecl 0x1d03fd0 </usr/include/_G_config.h:46:20, /usr/include/stdio.h:391:53> __arg '__va
57 |-FormatAttr 0x1d04500 <line:392:32, col:60> printf 3 0
58 `-NoThrowAttr 0x1d04550 </usr/include/sys/cdefs.h:57:37>
59 中间进行了省略
60 |-ParmVarDecl 0x1d05130 <<invalid sloc>> 'const char *restrict'
61 `-FormatAttr 0x1d05190 <col:12> scanf 1 2
62 FunctionDecl 0x1d051e0 prev 0x1d05090 <col:1, col:55> scanf 'int (const char *restrict, ...)' ext
63 |-ParmVarDecl 0x1d05010 <col:19, col:42> __format 'const char *restrict'
64 `-FormatAttr 0x1d052b0 <col:12> scanf 1 2
65 FunctionDecl 0x1d05470 <line:433:12> sscanf 'int (const char *restrict, const char *restrict, ..
66 |-ParmVarDecl 0x1d05510 <<invalid sloc>> 'const char *restrict'
67 |-ParmVarDecl 0x1d05570 <<invalid sloc>> 'const char *restrict'
68 `-FormatAttr 0x1d055e0 <col:12> scanf 2 3
```

clang 编译器的终端使用用颜色区分语法树。我们可以用图形显示如上：

LLVM IR 有 3 种表示形式：text: 便于阅读的文本格式，类似于汇编语言，拓展名.ll，$clang - S - emit - llvm main.m$ memory: bitcode: .bc clang -c -emit-llvm main.m

本次实验以 text 形式编译查看。编译器同时会使用语法树和符号表中信息来检查源程序是否与语言定义语义一致，进行类型检查等。在这一步，编译器同时进行了语法、词法相关的错误检查。若 mian.c 源程序中存在语法错误，clang 在输出语法树的同时也会报错。

## 3.3 中间代码生成与优化

具体命令如下：

图 3.1: clang 编译器语法树

```
1   gcc -fdump-tree-all-graph -fdump-rtl-all-graph main.c
2   dot -Tpng a-main.c.015t.cfg.dot > a-main.c.015t.cfg.dot.png
3   dot -Tpng a-main.c.244t.optimized.dot > a-main.c.244t.optimized.dot.png
```

在经过语法分析之后，编译器会生成一个明确的类机器语言的中间表示。实验中笔者通过-fdump-tree-all-graph 和-fdump-rtl-all-graph 两个 gcc flag 获得中间代码生成的多阶段的输出。由于生成了较多阶段的.dot 文件。笔者选择了部分生成的.dot 文件，利用 graphviz 进行可视化。可视化后我们可以看到程序的控制流图（CFG）以及各阶段处理中 CFG 的变化。

如图下图所示，`a-main.c.015t.cfg.dot`和`a-main.c.244t.optimized.dot`这两个文件分别是程序的控制流图和优化后的控制流图。我们不难发现，在优化后的程序中，程序的变量和相关罗杰逻辑被改变了，相对未优化的版本，优化的版本流程更加简短，减少了执行的时间。

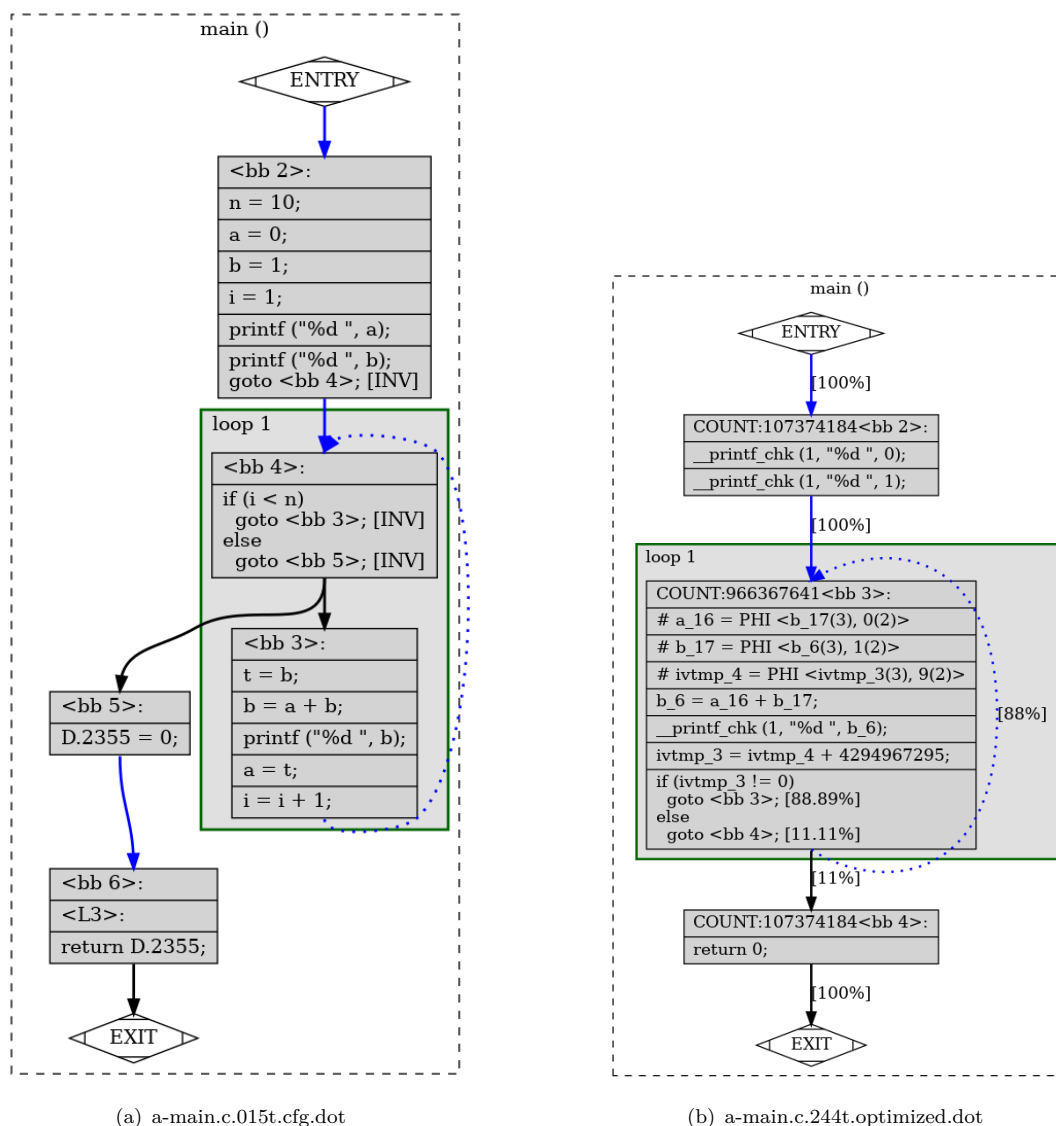(a) a-main.c.015t.cfg.dot　　　　　　　　　　　　　(b) a-main.c.244t.optimized.dot

图 3.2: GCC 中间结果的可视化

　　LLVM 的中间结果是一种类似汇编、却又不受限于平台的特殊语言 LLVM IR。笔者通过`-emit-llvm`参数来获取它。

```
1    clang -S -emit-llvm -O2 ./main.c -o main.O2.ll
2    clang -S -emit-llvm -Xclang -disable-O0-optnone ./main.c -o main.ll
3
4    opt -S main.ll -O1 -print-before-all > fin.log.ll 2>&1
```

　　本次实验通过`-O2`和`-Xclang -disable-O0-optnone`两个选项，输出流阶乘在 O2 优化和无优化选项下的 LLVM IR 代码，由于笔者有意使用了常量N作为斐波那契数列的长度，因此在 O2 优化下，编译器会在编译时就计算出数列的值，之后在运行时程序就可以直接将其输出。其中 clang -S -emit-llvm -O2 ./main.c -o main.O2.ll 为 O2 优化级别，clang -S -emit-llvm -Xclang -disable-O0-optnone ./main.c -o main.ll 为禁用优化级别，同时使用 opt -S main.ll -O1 -print-before-all > fin.log.ll 2>1 来输出优化

结果，编译器将会根据优化的结果生成汇编代码。

# 4　汇编器做了什么？

从本质上讲，汇编器也是编译器的一种，只是汇编器处理的"高级语言"是汇编语言，输出的是机器语言二进制形式。查阅资料发现，汇编器的构造和编译器相似，都需要进行词法分析、语法分析、语义处理、符号表管理和代码生成（机器代码）等阶段。

对于编译器来说，代码生成阶段只需要将解析的语法树映射到汇编语言子模块，而对于汇编器，将解析出的指令简洁的映射到正确机器代码相对比较复杂。由于本汇编器处理的输入文件为编译器生成的汇编文件，因此汇编器不需要考虑源文件会产生错误，因此它的语法分析的目的是识别出输入语言的语法结构并进行解析引导机器代码生成，因此不需要进行语法检查。汇编过程实际上把汇编语言程序代码翻译成目标机器指令的过程。其最终生成的是可重定位的机器代码。这一步一般被视为编译过程的"后端"。可以使用 gcc 命令 `gcc main.c -c -o main.o` 将汇编代码汇编为二进制文件。

通过 readelf 来解析这个文件的相关信息，如下图所示。

```
1   ELF Header:
2  Magic:   7f 45 4c 46 02 01 01 00 00 00 00 00 00 00 00 00
3  Class:                             ELF64
4  Data:                              2's complement, little endian
5  Version:                           1 (current)
6  OS/ABI:                            UNIX - System V
7  ABI Version:                       0
8  Type:                              REL (Relocatable file)
9  Machine:                           Advanced Micro Devices X86-64
10 Version:                           0x1
11 Entry point address:               0x0
12 Start of program headers:          0 (bytes into file)
13 Start of section headers:          832 (bytes into file)
14 Flags:                             0x0
15 Size of this header:               64 (bytes)
16 Size of program headers:           0 (bytes)
17 Number of program headers:         0
18 Size of section headers:           64 (bytes)
19 Number of section headers:         13
20 Section header string table index: 12
21
22 ction Headers:
23 [Nr] Name              Type            Address           Offset
24      Size              EntSize         Flags  Link  Info  Align
25 [ 0]                   NULL            0000000000000000  00000000
26      0000000000000000  0000000000000000           0     0     0
27 [ 1] .text             PROGBITS        0000000000000000  00000040
```

```
28       0000000000000066  0000000000000000  AX       0     0     1
29 [ 2] .rela.text         RELA              0000000000000000  00000260
30       0000000000000060  0000000000000018  I       10     1     8
31 [ 3] .data              PROGBITS          0000000000000000  000000a6
32       0000000000000000  0000000000000000  WA       0     0     1
33 [ 4] .bss               NOBITS            0000000000000000  000000a6
34       0000000000000000  0000000000000000  WA       0     0     1
35 [ 5] .rodata            PROGBITS          0000000000000000  000000a6
36       000000000000000a  0000000000000000  A        0     0     1
37 [ 6] .comment           PROGBITS          0000000000000000  000000b0
38       000000000000002e  0000000000000001  MS       0     0     1
39 [ 7] .note.GNU-stack    PROGBITS          0000000000000000  000000de
40       0000000000000000  0000000000000000           0     0     1
41 [ 8] .eh_frame          PROGBITS          0000000000000000  000000e0
42       0000000000000038  0000000000000000  A        0     0     8
43 [ 9] .rela.eh_frame     RELA              0000000000000000  000002c0
44       0000000000000018  0000000000000018  I       10     8     8
45 [10] .symtab            SYMTAB            0000000000000000  00000118
46       0000000000000120  0000000000000018          11     9     8
47 [11] .strtab            STRTAB            0000000000000000  00000238
48       0000000000000023  0000000000000000           0     0     1
49 [12] .shstrtab          STRTAB            0000000000000000  000002d8
50       0000000000000061  0000000000000000           0     0     1
51 y to Flags:
52 W (write), A (alloc), X (execute), M (merge), S (strings), I (info),
53 L (link order), O (extra OS processing required), G (group), T (TLS),
54 C (compressed), x (unknown), o (OS specific), E (exclude),
55 l (large), p (processor specific)
56
57 ere are no section groups in this file.
58
59 ere are no program headers in this file.
60
61 location section '.rela.text' at offset 0x260 contains 4 entries:
62 Offset          Info            Type            Sym. Value      Sym. Name + Addend
63 0000000010   00050000000a R_X86_64_32        0000000000000000 .rodata + 0
64 000000001a   000a00000002 R_X86_64_PC32      0000000000000000 __isoc99_scanf - 4
65 0000000051   00050000000a R_X86_64_32        0000000000000000 .rodata + 3
66 000000005b   000b00000002 R_X86_64_PC32      0000000000000000 printf - 4
67
68 location section '.rela.eh_frame' at offset 0x2c0 contains 1 entries:
69 Offset          Info            Type            Sym. Value      Sym. Name + Addend
```

```
70  0000000020   000200000002 R_X86_64_PC32      0000000000000000 .text + 0

71

72  e decoding of unwind sections for machine type Advanced Micro Devices X86-64 is not currently su

73

74  mbol table '.symtab' contains 12 entries:
75   Num:    Value          Size Type    Bind    Vis      Ndx Name
76     0: 0000000000000000     0 NOTYPE  LOCAL   DEFAULT  UND
77     1: 0000000000000000     0 FILE    LOCAL   DEFAULT  ABS main.c
78     2: 0000000000000000     0 SECTION LOCAL   DEFAULT    1
79     3: 0000000000000000     0 SECTION LOCAL   DEFAULT    3
80     4: 0000000000000000     0 SECTION LOCAL   DEFAULT    4
81     5: 0000000000000000     0 SECTION LOCAL   DEFAULT    5
82     6: 0000000000000000     0 SECTION LOCAL   DEFAULT    7
83     7: 0000000000000000     0 SECTION LOCAL   DEFAULT    8
84     8: 0000000000000000     0 SECTION LOCAL   DEFAULT    6
85     9: 0000000000000000   102 FUNC    GLOBAL  DEFAULT    1 main
86    10: 0000000000000000     0 NOTYPE  GLOBAL  DEFAULT  UND __isoc99_scanf
87    11: 0000000000000000     0 NOTYPE  GLOBAL  DEFAULT  UND printf

88

89   version information found in this file.

90
```

# 5   链接器做了什么？

由汇编程序生成的目标文件不能够直接执行。大型程序经常被分成多个部分进行编译，因此，可重定位的机器代码有必要和其他可重定位的目标文件以及库文件链接到一起，最终形成真正在机器上运行的代码。进而连接器对该机器代码进行执行生成可执行文件。

使用gcc main.o -o main命令，链接器根据汇编器中的信息将可重定位文件组合在一起，由于汇编器产生的可重定位文件中函数和其他标识符仍然有其名称定义，并不绑定到任何特定的地址，因此在将程序编译为可重定位 ELF 文件后，链接器生成可以链接可重定位文件会与库文件为一个整体，最生成可执行 ELF 文件。

# 6   LLVM IR 编程

## 6.1   小组分工

小组分工为：熊宇轩同学负责变量声明、赋值、运算符、数组等 SysY 特性，王天鹏同学负责循环、分支和函数调用等 SysY 特性，由于小组成员目前对编译相关知识了解还不多，后期相关任务分配可能会有调整。

但本次的 LLVM IR 编程作业任务较少，故使用王天鹏编写的 C 语言阶乘代码，由熊宇轩将其翻译成 LLVM IR 代码，此外，两人都参与了代码的调试和修饰工作。

## 6.2 C 语言阶乘源码

C 语言的阶乘程序源码如下所示：

```c
#include <stdio.h>

int factorial(int n)
{
    int i = 2, f = 1;
    if (n < 0)
        return -1;
    while(i <= n)
    {
        f = f*i;
        i = i+1;
    }
    return f;
}

int main()
{
    int n;
    scanf("%d", &n);
    n = factorial(n);
    printf("%d", n);
    return 0;
}
```

## 6.3 LLVM IR 代码

LLVM IR 的阶乘程序源码如下所示：

```llvm
; Function Attrs: noinline nounwind uwtable
define dso_local i32 @factorial(i32 noundef %0) #0 {
  ; ***BEGIN OF STUDENT CODE***
  %2 = alloca i32, align 4 ; i
  %3 = alloca i32, align 4 ; f
  %4 = alloca i32, align 4
  ; ↑没有用，但也删不掉，问就是非得按顺序来，
  ; 不然就'%5' is expected to be numbered '%4'
  ; 编译器级别的"这句没用，但是别改，删了跑不了"
  ; 重构时都是先用英语命名标识符，然后从上到下一个个分配数字
  ; 后来才发现label不一定要用数字。。。
  store i32 2, i32* %2, align 4 ; i = 2
  store i32 1, i32* %3, align 4 ; f = 1
```

```llvm
14
15    %5 = icmp slt i32 %0, 0 ; n < 0
16    br i1 %5, label %6, label %7
17
18  6:
19    store i32 −1, i32* %3, align 4 ; ret = −1
20    br label %15
21
22  7:
23    %8 = load i32, i32* %2, align 4 ; i
24    %9 = icmp sle i32 %8, %0 ; i <= n
25    ; ^ 这个i<=n之前写反了导致程序跑不了, 瞪了好久, 不知道怎么调试LLVM IR
26    ;   后来想起用Cutter/Ghidra反编译成C代码, 一眼发现。。。
27    br i1 %9, label %10, label %15
28
29  10:
30    %11 = load i32, i32* %2, align 4 ; i
31    %12 = load i32, i32* %3, align 4 ; f
32    %13 = mul nsw i32 %11, %12  ; f*i
33    store i32 %13, i32* %3, align 4 ; f = f*i
34    %14 = add nsw i32 %11, 1 ; i+1
35    store i32 %14, i32* %2, align 4 ; i = i+1
36    br label %7
37
38  15:
39    %16 = load i32, i32* %3, align 4 ; ret
40    ret i32 %16
41    ; ***END OF STUDENT CODE***
42  }
43
44  ; Function Attrs: noinline nounwind uwtable
45  define dso_local i32 @main() #0 {
46    ; ***BEGIN OF STUDENT CODE***
47    %1 = alloca i32, align 4 ; n
48    %2 = alloca i32, align 4 ; 13
49
50    %3 = call i32 (i8*, ...) @__isoc99_scanf(i8* noundef getelementptr inbounds ([3 x
          i8], [3 x i8]* @.str, i64 0, i64 0), i32* noundef %1)
51
52    %4 = load i32, i32* %1, align 4
53    %5 = call i32 @factorial(i32 noundef %4); 13 = factorial(n)
54
55    %6 = call i32 (i8*, ...) @printf(i8* noundef getelementptr inbounds ([3 x i8], [3 x
          i8]* @.str, i64 0, i64 0), i32 noundef %5)
56
57    ret i32 0
58
59    ; ***END OF STUDENT CODE***
60  }
```

　　上文中的代码仅包含了函数实现的部分，其他外围 LLVM IR 的代码则是照搬 clang 编译 C 语言版阶乘程序时所生成结果，完整的手工翻译的及 clang 生成的 LLVM IR 代码GitHub Gist找到。

　　通过`llc ./mytest.ll -filetype=obj -o test.o`和`clang ./test.o -static -o ./test`，我们成功将以上 LLVM IR 代码编译成了目标程序，经过测试，其行为与直接编译 C 语言版本的源码一致。