



南开大学
Nankai University

南 开 大 学
网 络 空 间 安 全 学 院
编译原理实验报告

第一次作业

王金晨 1911473

年级：2019 级

专业：信息安全

指导教师：王刚

2021 年 9 月 21 日

摘要

本文主要对我们经常使用的编译器 gcc 与 llvm 进行研究，来探究一个编译器在预处理阶段是如何对引用文件、命名空间、宏定义进行处理；在编译阶段如何进行词法分析、语法分析、语义分析、中间代码生成以及机器无关优化；以及汇编器与链接器加载器的工作探究。此外，还对于不同的优化级别在编译阶段中的各过程的中间结果的区别进行了部分探究。

关键字：预处理, 词法分析, 语法分析, 中间代码生成, 优化

目录

一、 概述	1
(一) 预处理	1
(二) 编译器	4
1. 词法分析	5
2. 语法分析	6
3. IR 生成	8
4. 机器无关优化	9
5. 汇编代码生成	10
(三) 汇编器	12
(四) 链接器加载器	13
二、 总结	15

一、 概述

(一) 预处理

在这里我们对编译器如何处理引用文件、静态变量声明定义与宏定义进行探究，推测静态定义与宏定义对文件引用部分无影响，宏定义的处理方式是对宏的替换，代码如下：

对照代码

```
1  #include<iostream>
2  using namespace std;
3  int main()
4  {
5      int a, b, i, n, t;
6      a=0;b=1;i=1;
7      cin>>n;
8      cout<<a<<endl;
9      cout<<b<<endl;
10     while(i<n)
11     {
12         t=b;
13         b = a + b;
14         cout<<b<<endl;
15         a=t;
16         i=i+1;
17     }
18     return 0;
19 }
```

静态定义

```
1  #include<iostream>
2  using namespace std;
3  static int a = 0;
4  static int b = 1;
5  static int t = 0;
6  int main()
7  {
8      int i, n;
9      cin>>n;
10     cout<<a<<endl;
11     cout<<b<<endl;
12     while(i<n)
13     {
14         t=b;
15         b = a + b;
16         cout<<b<<endl;
17         a=t;
18         i=i+1;
19     }
```

```
20     return 0;
21 }
```

宏定义

```
1 #include<iostream>
2 using namespace std;
3 #define zero int(0)
4 #define one int(1)
5 int main()
6 {
7     int a, b, i, n, t;
8     a=zero;b=one;i=one;
9     cin>>n;
10    cout<<a<<endl;
11    cout<<b<<endl;
12    while(i<n)
13    {
14        t=b;
15        b = a + b;
16        cout<<b<<endl;
17        a=t;
18        i=i+one;
19    }
20    return 0;
21 }
```

使用如下命令进行预处理:

```
gcc -E test.cpp > test.i
```

由于预处理结果有几万行, 在这里仅显示对照代码预处理后的几十行结果以及经过处理后的 main 函数代码:

pre

```
1 # 1 "test1.cpp"
2 # 1 "<built-in>"
3 # 1 "<command-line>"
4 # 1 "/usr/include/stdc-predef.h" 1 3 4
5 # 1 "<command-line>" 2
6 # 1 "test1.cpp"
7 # 1 "/usr/include/c++/7/iostream" 1 3
8 # 36 "/usr/include/c++/7/iostream" 3
9 # 37 "/usr/include/c++/7/iostream" 3
10 # 1 "/usr/include/x86_64-linux-gnu/c++/7/bits/c++config.h" 1 3
11 # 229 "/usr/include/x86_64-linux-gnu/c++/7/bits/c++config.h" 3
12 # 229 "/usr/include/x86_64-linux-gnu/c++/7/bits/c++config.h" 3
13 namespace std
14 {
15     typedef long unsigned int size_t;
```

```
16 typedef long int ptrdiff_t;
17 typedef decltype(nullptr) nullptr_t;
18 }
19 ...
20 # 2 "test1.cpp"
21 using namespace std;
22 int main()
23 {
24     int a, b, i, n, t;
25     a=0;b=1;i=1;
26     cin>>n;
27     cout<<a<<endl;
28     cout<<b<<endl;
29     while(i<n)
30     {
31         t=b;
32         b = a + b;
33         cout<<b<<endl;
34         a=t;
35         i=i+1;
36     }
37     return 0;
38 }
```

可以看到，gcc 将该 cpp 文件包含的文件编译进来，并对 namespace std 中的定义写入到本文件中。由于文件内容太多，所以我们使用 diff 命令来查看文件间的不同以便快速找到编译器处理 static 与宏命令的地方：其中，test1 是对照文件，test2 是添加了 static 定义的文件，test3 是使用了宏定义的文件。

diff test1 test2

```
1 1c1
2 < # 1 "test1.cpp"
3 ---
4 > # 1 "test2.cpp"
5 6c6
6 < # 1 "test1.cpp"
7 ---
8 > # 1 "test2.cpp"
9 28152c28152
10 < # 2 "test1.cpp" 2
11 ---
12 > # 2 "test2.cpp" 2
13 28154c28154
14 < # 2 "test1.cpp"
15 ---
16 > # 2 "test2.cpp"
17 28155a28156,28158
```

```

18 > static int a = 0;
19 > static int b = 1;
20 > static int t = 0;
21 28158,28159c28161
22 <     int a, b, i, n, t;
23 <     a=0;b=1;i=1;
24 ____
25 >     int i, n;

```

diff test1 test3

```

1 1c1
2 < # 1 "test1.cpp"
3 ____
4 > # 1 "test3.cpp"
5 6c6
6 < # 1 "test1.cpp"
7 ____
8 > # 1 "test3.cpp"
9 28152c28152
10 < # 2 "test1.cpp" 2
11 ____
12 > # 2 "test3.cpp" 2
13 28154c28154
14 < # 2 "test1.cpp"
15 ____
16 > # 2 "test3.cpp"
17 28155a28156,28157
18 >
19 >
20 28159c28161
21 <     a=0;b=1;i=1;
22 ____
23 >     a=int(0);b=int(1);i=int(1);
24 28169c28171
25 <     i=i+1;
26 ____
27 >     i=i+int(1);

```

可见，对于 static 定义，并没有进行其他处理，只是由于在编写代码时将 static 定义放在 main 函数开始之前导致的位置不同。对于宏的处理是把使用的宏直接进行文本意义上的替换，并删除宏定义语句。

(二) 编译器

在这一小节我探究了五个方面：词法分析、语法分析、中间代码生成（以下简称 IR 生成）、机器无关优化、目标代码生成（汇编代码）。在词法分析上探究了如何使用 gcc 与 llvm 进行 tokens 的扫描，在语法分析上主要探究如何使用上述二编译器生成 AST(抽象语法树)，IR 生成方面主要了解了 IR 的含义以及 CFG 对 IR 的反应，机器无关优化方面主要探究了 llvm 对 IR 的优化，

以及开启不同优化力度后 IR 的变化，目标代码生成则主要探究了 llvm 是如何生成汇编代码的。二编译器还有一个重要阶段是语义分析，该阶段作用为根据 AST 进行正确性上的检查，探究此部分需要阅读源代码，精力不足，不做探究。

1. 词法分析

使用命令：clang -E -Xclang -dump-tokens test.cpp

由于实际上是对预处理后的文件进行词法分析，所以会分析大量的 iostream 等库中的内容，导致词法分析结果过长，这里给出从分析 using name std 这句话开始的部分结果，如下：

```

1 using 'using'      [StartOfLine]  Loc=<test.cpp:2:1>
2 namespace 'namespace' [LeadingSpace] Loc=<test.cpp:2:7>
3 identifier 'std'      [LeadingSpace] Loc=<test.cpp:2:17>
4 semi ';'            Loc=<test.cpp:2:20>
5 int 'int'           [StartOfLine]  Loc=<test.cpp:3:1>
6 identifier 'main'     [LeadingSpace] Loc=<test.cpp:3:5>
7 l_paren '('         Loc=<test.cpp:3:9>
8 r_paren ')'         Loc=<test.cpp:3:10>
9 l_brace '{'         [StartOfLine]  Loc=<test.cpp:4:1>
10 int 'int'           [StartOfLine] [LeadingSpace] Loc=<test.cpp:5:5>
11 identifier 'a'       [LeadingSpace] Loc=<test.cpp:5:9>
12 comma ','           Loc=<test.cpp:5:10>
13 identifier 'b'       [LeadingSpace] Loc=<test.cpp:5:12>
14
15 equal '='           Loc=<test.cpp:6:10>
16 numeric_constant '1' Loc=<test.cpp:6:11>
17 semi ';'           Loc=<test.cpp:6:12>
18 identifier 'i'       Loc=<test.cpp:6:13>
19 equal '='           Loc=<test.cpp:6:14>
20 numeric_constant '1' Loc=<test.cpp:6:15>
21 semi ';'           Loc=<test.cpp:6:16>
22 identifier 'cin'     [StartOfLine] [LeadingSpace] Loc=<test.cpp:7:5>
23 greatergreater '>>' Loc=<test.cpp:7:8>
24 identifier 'n'       Loc=<test.cpp:7:10>
25 semi ';'           Loc=<test.cpp:7:11>
26 identifier 'cout'    [StartOfLine] [LeadingSpace] Loc=<test.cpp:8:5>
27 lessless '<<'       Loc=<test.cpp:8:9>
28 identifier 'a'       Loc=<test.cpp:8:11>
29 lessless '<<'       Loc=<test.cpp:8:12>
30 identifier 'endl'     Loc=<test.cpp:8:14>
31 semi ';'           Loc=<test.cpp:8:18>
32 ..
33 semi ';'           Loc=<test.cpp:9:18>
34 while 'while'       [StartOfLine] [LeadingSpace] Loc=<test.cpp:10:5>
35 l_paren '('         Loc=<test.cpp:10:10>
36 identifier 'i'       Loc=<test.cpp:10:11>
37 less '<'           Loc=<test.cpp:10:12>
38 identifier 'n'       Loc=<test.cpp:10:13>
39

```

```

40 plus '+' [LeadingSpace] Loc=<test.cpp:13:9>
41 identifier 'b' [LeadingSpace] Loc=<test.cpp:13:11>
42 semi ';' Loc=<test.cpp:13:12>
43 identifier 'cout' [StartOfLine] [LeadingSpace] Loc=<test.cpp:14:3>
44 lessless '<<' Loc=<test.cpp:14:7>
45 identifier 'b' Loc=<test.cpp:14:9>
46 lessless '<<' Loc=<test.cpp:14:10>
47 identifier 'endl' Loc=<test.cpp:14:12>
48 semi ';' Loc=<test.cpp:14:16>
49
50 numeric_constant '1' Loc=<test.cpp:16:7>
51 semi ';' Loc=<test.cpp:16:8>
52 r_brace '}' [StartOfLine] [LeadingSpace] Loc=<test.cpp:17:5>
53 return 'return' [StartOfLine] [LeadingSpace] Loc=<test.cpp:18:5>
54 numeric_constant '0' [LeadingSpace] Loc=<test.cpp:18:12>
55 semi ';' Loc=<test.cpp:18:13>
56 r_brace '}' [StartOfLine] Loc=<test.cpp:19:1>
57 eof '' Loc=<test.cpp:21:1>

```

可见。词法分析按照顺序进行扫描，为什么是从 using 开始呢？因为在实操中，发现这条命令输出了数千行。而 using 之前的部分都是对于预编译结果的解读，而从 using 还是才是对程序正文的解读，可以看到，词法分析首先忽略了所有的空格换行等无含义字符，随后将每个关键词（如 using、std）视为一个 token，那么在相应的 CFG(上下文无关文法) 中应该有这些终结符的存在，同时，他将每一个变量名字作为一个 token，甚至每个符号，可见这些都是 CFG 中的终结符。

2. 语法分析

主要探究 gcc 与 llvm 生成的语法树

对于 gcc 使用命令：gcc -fdump-tree-original-raw test.cpp

打开生成的文本文件，发现接近十几万行，而其中按照各个函数为划分构建的语法分析树，这里给出几个有代表性的节点：

```

1 ;; Function std::bad_exception::bad_exception() (null)
2 ;; enabled by -tree-original
3
4 @1      must_not_throw_expr type: @2      line: 49      body: @3
5 @2      void_type          name: @4      algn: 8
6 @3      statement_list    0 : @5      1 : @6
7 @4      type_decl         name: @7      type: @2      srcp: <built-in>:0
8                               note: artificial
9 @5      cleanup_point_expr type: @2      op 0: @8
10 @6      bind_expr         type: @2      body: @9
11 @7      identifier_node   strg: void    lngt: 4
12 @8      expr_stmt         type: @2      expr: @10
13 @9      statement_list    0 : @11     1 : @12
14 @10     modify_expr       type: @2      op 0: @13     op 1: @14
15 @11     cleanup_point_expr type: @2      op 0: @15
16 @12     try_catch_expr    type: @2

```



```

17 @13 indirect_ref type: @16 op 0: @17
18 @14 constructor lngt: 0
19 @15 expr_stmt type: @2 line: 49 expr: @18
20 @16 record_type size: @19 algn: 64 bfld: @20
21
22 @24 addr_expr type: @35 op 0: @36
23 @29 function_decl name: @41 type: @46 scpe: @20
24 srcp: exception:49 note: member
25 accs: pub note: constructor
26 args: @22 link: extern
27 @65 identifier_node strg: exception lngt: 9
28 @66 parm_decl name: @31 type: @67 scpe: @52
29 srcp: exception.h:63 note: artificial
30 argt: @67 size: @19 algn: 64
31 used: 0
32 @67 pointer_type qual: c unql: @35 size: @19
33 algn: 64 ptd: @45
34 @68 nop_expr type: @48 op 0: @22
35 @69 type_decl name: @74 type: @56 srcp: <built-in>:0
36 note: artificial
37 @70 function_type size: @61 algn: 8 retn: @75
38 @71 method_type unql: @76 size: @61 algn: 8
39 clas: @45 retn: @2 prms: @77
40 @75 integer_type name: @78 size: @79 algn: 32
41 prec: 32 sign: signed min: @80
42 max: @81
43 @76 method_type size: @61 algn: 8 clas: @45
44 retn: @2 prms: @77
45 @77 tree_list valu: @35 chan: @82
46 @79 integer_cst type: @25 int: 32
47 @82 tree_list valu: @84 chan: @73
48 @83 identifier_node strg: int lngt: 3
49 @84 reference_type size: @19 algn: 64 refd: @85
50 @85 record_type qual: c name: @57 unql: @45
51 size: @19 algn: 64 vfld: @27
52 tag: struct flds: @27 binf: @59

```

经过分析，每一行为一个节点 @ 加数字即表示一个节点，那么这到底有什么含义呢？经查阅资料，gcc 的 AST 的树节点结构如下：

```

1 union tree_node
2 {
3     ....
4 }

```

节点内部内容省略，也就是说，这里的每个节点是一个 union 结构，类似于 C++ 中的类，但只占用一个成员大小的内存，也就是共用。

对于 llvm 使用命令：clang -E -Xclang -ast-dump test.cpp

生成的部分语法树如下：

```

| | | | -ImplicitCastExpr 0x55e37b35b808 <col:11> 'int' <LValueToRValue>
| | | | -DeclRefExpr 0x55e37b35b7c8 <col:11> 'int' lvalue Var 0x55e37b348c40 'b' 'int'
| | | | -CXXOperatorCallExpr 0x55e37b35d028 <line:14:3, col:12> 'std::basic_ostream<char> >::_
| | | | _ostream_type': 'std::basic_ostream<char>' lvalue
| | | | -ImplicitCastExpr 0x55e37b35d010 <col:10> 'std::basic_ostream<char, std::char_traits<char> >::_ostream_type
| | | | e &(*) (std::basic_ostream<char, std::char_traits<char> >::_ostream_type &(*) (std::basic_ostream<char, std::char_traits<char
| | | | s<char> >::_ostream_type &*))' <FunctionToPointerDecay>
| | | | -DeclRefExpr 0x55e37b35cfe8 <col:10> 'std::basic_ostream<char, std::char_traits<char> >::_ostream_type &
| | | | (std::basic_ostream<char, std::char_traits<char> >::_ostream_type &(*) (std::basic_ostream<char, std::char_traits<char
| | | | > >::_ostream_type &*))' lvalue CXXMethod 0x55e37b2a2b00 'operator<<' 'std::basic_ostream<char, std::char_traits<char>
| | | | >::_ostream_type &(std::basic_ostream<char, std::char_traits<char> >::_ostream_type &(*) (std::basic_ostream<char, s
| | | | td::char_traits<char> >::_ostream_type &*))'
| | | | -CXXOperatorCallExpr 0x55e37b35c648 <col:3, col:9> 'std::basic_ostream<char, std::char_traits<char> >::_os
| | | | tream_type': 'std::basic_ostream<char>' lvalue
| | | | -ImplicitCastExpr 0x55e37b35c630 <col:7> 'std::basic_ostream<char, std::char_traits<char> >::_ostream_ty
| | | | pe &(*) (int)' <FunctionToPointerDecay>
| | | | -DeclRefExpr 0x55e37b35c608 <col:7> 'std::basic_ostream<char, std::char_traits<char> >::_ostream_type
| | | | &(int)' lvalue CXXMethod 0x55e37b2a3b00 'operator<<' 'std::basic_ostream<char, std::char_traits<char> >::_ostream_ty
| | | | pe &(int)'
| | | | -DeclRefExpr 0x55e37b35b870 <col:3> 'std::ostream': 'std::basic_ostream<char>' lvalue Var 0x55e37b3484a8 '
| | | | cout' 'std::ostream': 'std::basic_ostream<char>'
| | | | -ImplicitCastExpr 0x55e37b35c5f0 <col:9> 'int' <LValueToRValue>
| | | | -DeclRefExpr 0x55e37b35b898 <col:9> 'int' lvalue Var 0x55e37b348c40 'b' 'int'
| | | | -ImplicitCastExpr 0x55e37b35cfd0 <col:12> 'basic_ostream<char, std::char_traits<char> > &(*) (basic_ostream<
| | | | char, std::char_traits<char> > &)' <FunctionToPointerDecay>
| | | | -DeclRefExpr 0x55e37b35cfa0 <col:12> 'basic_ostream<char, std::char_traits<char> > &(*) (basic_ostream<char,
| | | | std::char_traits<char> > &)' lvalue Function 0x55e37b2a73c0 'endl' 'basic_ostream<char, std::char_traits<char> > &(bas
| | | | ic_ostream<char, std::char_traits<char> > &)' (FunctionTemplate 0x55e37b286a20 'endl')
| | | | -BinaryOperator 0x55e37b35d0d8 <line:15:3, col:5> 'int' lvalue '='
| | | | -DeclRefExpr 0x55e37b35d070 <col:3> 'int' lvalue Var 0x55e37b348bc8 'a' 'int'
| | | | -ImplicitCastExpr 0x55e37b35d0c0 <col:5> 'int' <LValueToRValue>
| | | | -DeclRefExpr 0x55e37b35d098 <col:5> 'int' lvalue Var 0x55e37b348da8 't' 'int'
| | | | -BinaryOperator 0x55e37b35d1b0 <line:16:3, col:7> 'int' lvalue '='
| | | | -DeclRefExpr 0x55e37b35d100 <col:3> 'int' lvalue Var 0x55e37b348cb8 'i' 'int'
| | | | -BinaryOperator 0x55e37b35d188 <col:5, col:7> 'int' '+'
| | | | -ImplicitCastExpr 0x55e37b35d170 <col:5> 'int' <LValueToRValue>
| | | | -DeclRefExpr 0x55e37b35d128 <col:5> 'int' lvalue Var 0x55e37b348cb8 'i' 'int'

```

可以看出一定的树状结构，甚至标出了行号列号，但是由于语法树太长，这里看不到根节点。最顶部一般是 TranslationUnitDecl(一个 Cpp 文件以及那些 include 包括的文件称为翻译单元 (TranslationUnit))，TypedefDecl、CXXRecordDecl、CompoundStmt 等称为 AST node，比较常见的有 Stmt、Decl 和 Expr 等。[1]

3. IR 生成

llvm 的 IR 生成，使用命令：clang -S -emit-llvm test.cpp 生成 test.ll 文件除去一些声明与定义外，直接查看对应主函数的 IR，仅给出部分：

```

1 define i32 @main() #4 {
2   %1 = alloca i32, align 4
3   %2 = alloca i32, align 4
4   %3 = alloca i32, align 4
5
6   store i32 1, i32* %3, align 4
7   store i32 1, i32* %4, align 4
8   %7 = call dereferenceable(280) @"class.std::basic_istream"* @_ZNSirsERi(%"
   class.std::basic_istream"* @_ZSt3cin, i32* dereferenceable(4) %5)
9   %8 = load i32, i32* %2, align 4
10  %9 = call dereferenceable(272) @"class.std::basic_ostream"* @_ZNSolsEi(%"
   class.std::basic_ostream"* @_ZSt4cout, i32 %8)
11
12  %13 = call dereferenceable(272) @"class.std::basic_ostream"*
   @_ZNSolsEPFRSoS_E(%"class.std::basic_ostream"* %12, %"class.std::
   basic_ostream"* (%"class.std::basic_ostream"*)*
   @_ZSt4endlIcSt11char_traitsIcEERSt13basic_ostreamIT_T0_ES6_)
13  br label %14
14

```

```

15 ; <label>:14:                                ; preds = %18, %0
16   %15 = load i32, i32* %4, align 4
17   %16 = load i32, i32* %5, align 4
18   %17 = icmp slt i32 %15, %16
19   br i1 %17, label %18, label %29
20
21 ; <label>:18:                                ; preds = %14
22   %19 = load i32, i32* %3, align 4
23   store i32 %19, i32* %6, align 4
24   @_ZNSolsEPFRSoS_E(%"class.std::basic_ostream"* %24, %"class.std::
    basic_ostream"* (%"class.std::basic_ostream"*)*
    @_ZSt4endlIcSt11char_traitsIcEERSt13basic_ostreamIT_T0_ES6_)
25   %26 = load i32, i32* %6, align 4
26   store i32 %26, i32* %2, align 4
27   %27 = load i32, i32* %4, align 4
28   %28 = add nsw i32 %27, 1
29   store i32 %28, i32* %4, align 4
30   br label %14
31
32 ; <label>:29:                                ; preds = %14
33   ret i32 0
34 }

```

分析可知，% 开头的为局部变量，@ 开头的为全局变量，与我们熟知的汇编类似，load 从内存取，大部分指令按照英文意思就可以猜出，可以说是十分易于理解了。

4. 机器无关优化

使用命令：llc -print-before-all -print-after-all a.ll，部分结果如下：

```

1 main:                                         # @main
2   .cfi_startproc
3 # BB#0:
4   pushq   %rbp
5   .Ltmp3:
6   .cfi_def_cfa_offset 16
7   .Ltmp4:
8   .cfi_offset %rbp, -16
9   movq    %rsp, %rbp
10  .Ltmp5:
11  .cfi_def_cfa_register %rbp
12  subq     $32, %rsp
13  movabsq  $__ZSt3cin, %rdi
14  leaq     -20(%rbp), %rsi
15  movl     $0, -4(%rbp)
16  movl     $0, -8(%rbp)
17  movl     $1, -12(%rbp)
18  movl     $1, -16(%rbp)
19  callq    _ZNSirsERi

```

```

20     movabsq $__ZSt4cout, %rdi
21     movl    -8(%rbp), %esi
22     callq   __ZNSolsEi
23     movabsq $__ZSt4endlIcSt11char_traitsIcEERSt13basic_ostreamIT_T0_ES6_,
        %rsi
24     movq    %rax, %rdi
25     callq   __ZNSolsEPFRSoS_E
26     movabsq $__ZSt4cout, %rdi
27
28     movl    -16(%rbp), %eax
29     addl    $1, %eax
30     movl    %eax, -16(%rbp)
31     jmp     .LBB1_1
32 .LBB1_3:
33     xorl    %eax, %eax
34     addq    $32, %rsp
35     popq    %rbp
36     retq
37 .Lfunc_end1:
38     .size   main, .Lfunc_end1-main
39     .cfi_endproc

```

可以看到，优化的版本对于堆栈的使用更加高效，并且新增很多标签更加易读，另外，对于多数寄存器的操作，从操作整个 64 位寄存器优化为只操作寄存器的低 32 位，数据传输更快。

5. 汇编代码生成

由于使用的是 C++ 程序，所以使用 g++，g++ 是 gcc 的前端。命令如下：g++ -S test.i -o test.S 展示部分代码：

```

1     .file    "test.cpp"
2     .text
3     .section        .rodata
4     .type   __ZStL19piecewise_construct, @object
5     .size   __ZStL19piecewise_construct, 1
6     __ZStL19piecewise_construct:
7     .zero   1
8     .local  __ZStL8__ioinit
9     .comm   __ZStL8__ioinit,1,1
10    .text
11    .globl  main
12    .type   main, @function
13    main:
14    .LFB1493:
15        .cfi_startproc
16        pushq   %rbp
17        .cfi_def_cfa_offset 16
18        .cfi_offset 6, -16
19        movq    %rsp, %rbp

```

```

20      .cfi_def_cfa_register 6
21      subq    $32, %rsp
22      movq    %fs:40, %rax
23      movq    %rax, -8(%rbp)
24      xorl    %eax, %eax
25
26      movq    %rax, %rsi
27      movq    %rdx, %rdi
28      call    __ZNSolsEPFRSoS_E@PLT
29  .L3:
30      movl    -28(%rbp), %eax
31      cmpl    %eax, -16(%rbp)
32      jge     .L2
33      addl    $1, -16(%rbp)
34      jmp     .L3
35  .L2:
36      movl    $0, %eax
37      movq    -8(%rbp), %rcx
38      xorq    %fs:40, %rcx
39      je      .L5
40      call    ___stack_chk_fail@PLT
41  .L5:
42      leave
43      .cfi_def_cfa 7, 8
44      ret
45      .cfi_endproc

```

生成了 x86 汇编代码，可以看到，有在 main 之前的入口程序与 main 结束之后的收尾程序。再尝试生成 arm 格式汇编代码，同样只给出部分：

```

1      .arch armv7-a
2      .eabi_attribute 28, 1
3      .fpu vfpv3-d16
4      .eabi_attribute 20, 1
5      ...
6      .eabi_attribute 18, 4
7      .file    "test.cpp"
8      .local   __ZStL8__ioinit
9      ...
10     main:
11         .fnstart
12     .LFB1024:
13         @ args = 0, pretend = 0, frame = 24
14         @ frame_needed = 1, uses_anonymous_args = 0
15         push    {r7, lr}
16         .save   {r7, lr}
17         .pad    #24
18         sub     sp, sp, #24
19         setfp   r7, sp, #0

```

```

20      add     r7, sp, #0
21      movw   r3, #lower16: __stack_chk_guard
22      movt   r3, #upper16: __stack_chk_guard
23      ldr    r3, [r3]
24      str    r3, [r7, #20]
25      movs   r3, #0
26      ...
27
28      bl     __ZNSolsEi
29      mov    r3, r0
30      movw   r1, #lower16:
31      _ZSt4endlIcSt11char_traitsIcEERSt13basic_ostreamIT_T0_ES6_
32      movt   r1, #upper16:
33      _ZSt4endlIcSt11char_traitsIcEERSt13basic_ostreamIT_T0_ES6_
34      mov    r0, r3
35      bl     __ZNSolsEPFRSoS_E
36      ldr    r3, [r7, #16]
37      str    r3, [r7, #4]
38      ...
39      .fnend
40      .size   main, .-main
41      .align  2
42      .syntax unified
43      .thumb
44      .thumb_func
45      .type   __Z41__static_initialization_and_destruction_0ii, %function
46      __Z41__static_initialization_and_destruction_0ii:
47      .fnstart

```

但是在实验过程中,并不能由之前生成的 test.i 文件直接使用命令 arm-linux-gnueabi-g++ test.i -S -o test.S 生成汇编代码,需要先使用 arm-linux-gnueabi-g++ 进行预处理生成对应的 test.i 文件之后才可以,推测是二者对于 C++ 文件的预处理方式不同,根据报错分析可能是符号上或定义上了一些区别。

(三) 汇编器

使用 gcc 对 test.cpp(而不是 arm-linux-gnueabi-g++) 进行预处理与编译后,使用 gcc test.S -c -o test.o 生成二进制文件,这个时候的二进制文件还不是可执行的,因为他只是在源程序预处理之上所得到的二进制文件,本质上并不包含所引用的库的函数,所以也就没有办法正常执行,使用 objdump 反汇编,给出部分:

```

1  test.o:          文件格式 elf64-x86-64
2
3
4  Disassembly of section .text:
5
6  0000000000000000 <main>:
7      0:  55                push    %rbp
8      1:  48 89 e5          mov     %rsp,%rbp

```

9	4:	48 83 ec 20	sub	\$0x20,%rsp	
10	8:	64 48 8b 04 25 28 00	mov	%fs:0x28,%rax	
11	f:	00 00			
12	11:	48 89 45 f8	mov	%rax,-0x8(%rbp)	
13	15:	31 c0	xor	%eax,%eax	
14	17:	c7 45 e8 00 00 00 00	movl	\$0x0,-0x18(%rbp)	
15	1e:	c7 45 ec 01 00 00 00	movl	\$0x1,-0x14(%rbp)	
16	25:	c7 45 f0 01 00 00 00	movl	\$0x1,-0x10(%rbp)	
17	2c:	48 8d 45 e4	lea	-0x1c(%rbp),%rax	
18	30:	48 89 c6	mov	%rax,%rsi	
19	33:	48 8d 3d 00 00 00 00	lea	0x0(%rip),%rdi	# 3a <main+0x3a>
20	3a:	e8 00 00 00 00	callq	3f <main+0x3f>	
21	3f:	8b 45 e8	mov	-0x18(%rbp),%eax	
22	42:	89 c6	mov	%eax,%esi	
23	44:	48 8d 3d 00 00 00 00	lea	0x0(%rip),%rdi	# 4b <main+0x4b>
24	4b:	e8 00 00 00 00	callq	50 <main+0x50>	
25	50:	48 89 c2	mov	%rax,%rdx	
26	53:	48 8b 05 00 00 00 00	mov	0x0(%rip),%rax	# 5a <main+0x5a>
27	5a:	48 89 c6	mov	%rax,%rsi	
28	5d:	48 89 d7	mov	%rdx,%rdi	
29	60:	e8 00 00 00 00	callq	65 <main+0x65>	
30	65:	8b 45 ec	mov	-0x14(%rbp),%eax	
31	68:	89 c6	mov	%eax,%esi	
32	6a:	48 8d 3d 00 00 00 00	lea	0x0(%rip),%rdi	# 71 <main+0x71>
33	71:	e8 00 00 00 00	callq	76 <main+0x76>	
34	76:	48 89 c2	mov	%rax,%rdx	
35	79:	48 8b 05 00 00 00 00	mov	0x0(%rip),%rax	# 80

与原来的.S 文件有较大区别，多数体现在对于变量的命名以及一些声明定义等，原因猜测有 2：首先，编译器在将汇编语言转换成二进制时会进行一定程度的优化；其次，编译器在翻译时会导致一些信息的丢失，比如节等信息的设置等，会导致反汇编得到的结果不可能完全一样，因为这是信息不对等的。接下来开启优化尝试结果发现反汇编得到的代码一模一样，要么是在这一步不存在优化环节，要么是这个小程序没有优化空间了，至少在 o2 程度上没有优化空间

(四) 链接器加载器

进行链接：使用 `g++ test.o -o test`，而后使用 `objdump` 进行反汇编得到，给出部分：

1	test:	文件格式 elf64-x86-64			
2	Disassembly of section .init:				
3	00000000000007e0 <_init>:				
4	7e0:	48 83 ec 08	sub	\$0x8,%rsp	
5	7e4:	48 8b 05 fd 07 20 00	mov	0x2007fd(%rip),%rax	# 200fe8 <__gmon_start__>
6	7eb:	48 85 c0	test	%rax,%rax	
7	7ee:	74 02	je	7f2 <_init+0x12>	
8	7f0:	ff d0	callq	*%rax	
9	7f2:	48 83 c4 08	add	\$0x8,%rsp	
10	7f6:	c3	retq		

```

11 Disassembly of section .plt:
12 0000000000000800 <.plt>:
13 800: ff 35 82 07 20 00 pushq 0x200782(%rip) # 200f88 <
    _GLOBAL_OFFSET_TABLE_+0x8>
14 806: ff 25 84 07 20 00 jmpq *0x200784(%rip) # 200f90 <
    _GLOBAL_OFFSET_TABLE_+0x10>
15 80c: 0f 1f 40 00 nopl 0x0(%rax)
16
17 0000000000000810 <_ZNSirsERi@plt>:
18 810: ff 25 82 07 20 00 jmpq *0x200782(%rip) # 200f98 <
    _ZNSirsERi@GLIBCXX_3.4>
19 816: 68 00 00 00 00 pushq $0x0
20 81b: e9 e0 ff ff jmpq 800 <.plt>
21
22 0000000000000820 <__cxa_atexit@plt>:
23 820: ff 25 7a 07 20 00 jmpq *0x20077a(%rip) # 200fa0 <
    __cxa_atexit@GLIBC_2.2.5>
24 826: 68 01 00 00 00 pushq $0x1
25 82b: e9 d0 ff ff jmpq 800 <.plt>
26
27 ...
28
29 Disassembly of section .plt.got:
30
31 0000000000000870 <__cxa_finalize@plt>:
32 870: ff 25 52 07 20 00 jmpq *0x200752(%rip) # 200fc8 <
    __cxa_finalize@GLIBC_2.2.5>
33 876: 66 90 xchg %ax,%ax
34 Disassembly of section .text:
35 0000000000000880 <_start>:
36 880: 31 ed xor %ebp,%ebp
37 882: 49 89 d1 mov %rdx,%r9
38 885: 5e pop %rsi
39 ...
40 971: 0f 1f 80 00 00 00 00 nopl 0x0(%rax)
41 978: f3 c3 repz retq
42 97a: 66 0f 1f 44 00 00 nopw 0x0(%rax,%rax,1)
43
44
45 000000000000098a <main>:
46 98a: 55 push %rbp
47 98b: 48 89 e5 mov %rsp,%rbp
48 98e: 48 83 ec 20 sub $0x20,%rsp
49 992: 64 48 8b 04 25 28 00 mov %fs:0x28,%rax
50 999: 00 00
51 99b: 48 89 45 f8 mov %rax,-0x8(%rbp)
52 99f: 31 c0 xor %eax,%eax
53 9a1: c7 45 e8 00 00 00 00 movl $0x0,-0x18(%rbp)

```



```

54  9a8:  c7 45 ec 01 00 00 00    movl    $0x1,-0x14(%rbp)
55  9af:  c7 45 f0 01 00 00 00    movl    $0x1,-0x10(%rbp)
56  9b6:  48 8d 45 e4             lea      -0x1c(%rbp),%rax
57  ...
58  a4f:  8b 45 f4             mov     -0xc(%rbp),%eax
59  a52:  89 45 e8             mov     %eax,-0x18(%rbp)
60  a55:  83 45 f0 01          addl    $0x1,-0x10(%rbp)
61  a59:  eb ba             jmp     a15 <main+0x8b>
62  a5b:  b8 00 00 00 00       mov     $0x0,%eax
63  a60:  48 8b 4d f8          mov     -0x8(%rbp),%rcx
64  a64:  64 48 33 0c 25 28 00  xor     %fs:0x28,%rcx
65  a6b:  00 00
66  a6d:  74 05             je      a74 <main+0xea>
67  a6f:  e8 cc fd ff ff       callq   840 <__stack_chk_fail@plt>
68  a74:  c9             leaveq  %eax
69  a75:  c3             retq
70
71  Disassembly of section .fini:
72
73  0000000000000b54 <_fini>:
74  b54:  48 83 ec 08          sub     $0x8,%rsp
75  b58:  48 83 c4 08          add     $0x8,%rsp
76  b5c:  c3             retq

```

可以看到还是有较大不同的，test 文件的反汇编显示 init 程序，而后是 plt 表、start 程序，在 start 程序中为 main 函数分配了栈之后才正式执行 main 函数，并且在 main 函数之后会返回 start 函数并跳转到 fini 函数进行堆栈回收。同时，对于整个程序来说，多出了很多在源代码中没有使用到的函数，比如检查堆栈的函数等，这些是在链接加载过程中会给程序增加一些保护，比如数组越界等检查，来让程序有能力应对针对程序脆弱性的攻击，从而提高程序安全性。

二、 总结

预处理会将包含文件代码插入源文件，导致预处理结果很长，其主要工作就是插入引用的文件或库，并对宏做直接的替换。

编译阶段主要预处理的结果进行词法分析、语法分析、语义分析、优化、代码生成。词法分析采用逐个扫描识别 token，其中会对空格、Tab 键这种排版符号做忽视处理。语法分析主要生成抽象语法树，gcc 与 llvm 生成的语法树采用的节点定义并不相同，相比之下我认为 gcc 更为好读，且可以生成图片便于理解。优化方面可以看到 llvm 的 pass 在内存使用效率、寄存器使用效率、指令执行效率上进行了优化。在代码生成上可以看到由于平台架构不同，对于 x86 与 amd 架构需要不同的编译器来进行汇编代码的生成，原因在于二者汇编代码的不同。

汇编阶段对汇编代码进行处理，将其转为机器码，这个时候的二进制文件是不可执行的，其不具备库函数，也就没办法正常执行，同时也没有被指明在执行时应该被装载到哪里。在这个阶段是无法使用 gcc 对文件进行汇编的，原因在于使用的是 C++ 语言，gcc 无法找到 C++ 语言的库，因而使用 g++ 进行汇编，g++ 与 gcc 并无本质区别，g++ 可以看作对 gcc 的一个 C++ 上的封装，二者在预处理、编译阶段没有区别，只在生成二进制文件的时候产生区别，这是语言特性的原因。

链接阶段就没什么可说的了。

对于 gcc 与 llvm 二者在其他方面没有太大区别，除了灵活性上。llvm 采用统一的中间代码，使得其可以用于各个语言的编译，而不单单是用于 C 语言的 clang，其前端与后端可以自由选择，而不是向 gcc 一样对于不同架构需要开发不同版本。

NIKU

参考文献

- [1] liveforthing. Clang 之语法抽象语法树 AST. <https://www.cnblogs.com/zhangke007/p/4714245.html>, 2015.

NIKU