

C++编译链接过程

```
1 #include <Stdio.h>
2 int main(){
3     printf("hello world!\n");
4 }
```

```
1 gcc hello.c #编译
2 ./a.out #执行
```

上述gcc命令实际执行了四步操作：

1. 预处理 (Preprocessing)
2. 编译 (Compilation)
3. 汇编 (Assemble)
4. 链接 (Linking)

示例

```
1 |— test.c
2 |— inc
3     |— mymath.h
4     |— mymath.c
```

假设文件结构如上，test中使用mymath实现的函数

```
1 //test.c
2 #include <stdio.h>
3 #include "mymath.h"// 自定义头文件
4 int main(){
5     int a = 2;
6     int b = 3;
7     int sum = add(a, b);
8     printf("a=%d, b=%d, a+b=%d\n", a, b, sum);
9 }
```

```
1 // mymath.h
2 #ifndef MYMATH_H
3 #define MYMATH_H
4 int add(int a, int b);
5 int sum(int a, int b);
6 #endif
```

```
1 // mymath.c
2 int add(int a, int b){
3     return a+b;
4 }
5 int sub(int a, int b){
6     return a-b;
7 }
```

1. 预处理 (Preprocessing)

预处理用于将所有的**#include头文件以及宏定义替换成真正的内容**，预处理之后得到的仍然是文本文件，但文件体积会变大很多。gcc的预处理是**预处理器cpp**完成的，可通过如下命令对test.c进行预处理

```
1 | $ gcc -E -I./inc test.c -o test.i
```

或者直接调用cpp命令

```
1 | $ cpp test.c -I./inc -o test.i
```

上述命令中-E是让编译器在预处理之后就退出，不进行后续编译过程；-I指定头文件目录，这里指定的是我们自定义的头文件目录；-o指定输出文件名

2. 编译 (Compilation)

这里的编译不是指程序从源文件到二进制程序的全部过程，而是指将经过预处理之后的程序转换成特定汇编代码的过程

```
1 | $ gcc -S -I./inc test.c -o test.s
```

上述命令中-S让编译器在编译之后停止，不进行后续过程。编译过程完成后，**将生成程序的汇编代码test.s，这也是文本文件**，生成的指令如下

```
1 | .file "main.cpp"
2 | .text
3 | .section .text$_Z6printfPKcz,"x"
```

```
4      .linkonce discard
5      .globl  _Z6printfPKcz
6      .def    _Z6printfPKcz; .sc1    2; .type
32; .endef
7      .seh_proc  _Z6printfPKcz
8  _Z6printfPKcz:
9  .LFB9:
10     pushq    %rbp
11     .seh_pushreg    %rbp
12     pushq    %rbx
13     .seh_pushreg    %rbx
14     subq     $56, %rsp
15     .seh_stackalloc 56
16     leaq     48(%rsp), %rbp
17     .seh_setframe    %rbp, 48
18     .seh_endprologue
19     movq     %rcx, 32(%rbp)
20     movq     %rdx, 40(%rbp)
21     movq     %r8, 48(%rbp)
22     movq     %r9, 56(%rbp)
23     leaq     40(%rbp), %rax
24     movq     %rax, -16(%rbp)
25     movq     -16(%rbp), %rbx
26     movl     $1, %ecx
27     movq     __imp___acrt_iob_func(%rip), %rax
28     call     *%rax
29     movq     %rax, %rcx
30     movq     32(%rbp), %rax
31     movq     %rbx, %r8
32     movq     %rax, %rdx
```

```

33     call    __mingw_vfprintf
34     movl    %eax, -4(%rbp)
35     movl    -4(%rbp), %eax
36     addq    $56, %rsp
37     popq    %rbx
38     popq    %rbp
39     ret
40     .seh_endproc
41     .def    __main; .sc1    2; .type    32;
42     .endif
43     .section .rdata,"dr"
44     .LC0:
45     .ascii  "a=%d, b=%d, a+b=%d\12\0"
46     .text
47     .globl  main
48     .def    main; .sc1    2; .type    32;
49     .endif
50     .seh_proc  main
51     main:
52     .LFB45:
53     pushq   %rbp
54     .seh_pushreg    %rbp
55     movq    %rsp, %rbp
56     .seh_setframe   %rbp, 0
57     subq    $48, %rsp
58     .seh_stackalloc 48
59     .seh_endprologue
60     call    __main
61     movl    $2, -4(%rbp)
62     movl    $3, -8(%rbp)

```

```
61     movl    -8(%rbp), %edx
62     movl    -4(%rbp), %eax
63     movl    %eax, %ecx
64     call    _Z3addii
65     movl    %eax, -12(%rbp)
66     movl    -12(%rbp), %ecx
67     movl    -8(%rbp), %edx
68     movl    -4(%rbp), %eax
69     movl    %ecx, %r9d
70     movl    %edx, %r8d
71     movl    %eax, %edx
72     leaq    .LC0(%rip), %rax
73     movq    %rax, %rcx
74     call    _Z6printfPKcz
75     movl    $0, %eax
76     addq    $48, %rsp
77     popq    %rbp
78     ret
79     .seh_endproc
80     .ident   "GCC: (x86_64-win32-seh-rev1,
Built by MinGW-w64 project) 12.2.0"
81     .def     __mingw_vfprintf;    .sc1    2;
.type    32; .endef
82     .def     _Z3addii;    .sc1    2; .type
32; .endef
```

3. 汇编 (Assemble)

汇编过程将上一步的汇编代码转换成机器码，产生的文件叫做目标文件，是二进制格式

```
1 | as test.s -o test.o
```

等价于

```
1 | gcc -c test.s -o test.o
```

这一步会为每一个源文件产生一个目标文件，因此mymath.c也需要产生一个mymath.o文件

4. 链接 (Linking)

```
1 | $ ld -o test.out test.o inc/mymath.o  
...libraries...
```

链接的详细过程如下

合并段

在elf文件中字节对齐是以4字节对齐的，在可执行程序中是以页的方式对齐的（一个页的大小为4k），因此如果我们在链接时将各个.o文件各个段单独的加载到可执行文件中，将会非常浪费空间：

因此我们需要合并段，调整段偏移，把每个.o文件的.text段合并在一起.data段合并在一起，这样在生成的可执行文件中，各个段都只有一个，如下图，由于在链接时只需要加载代码段(.text段)和数据段（.data段和.bss段）。因此合并段之后，在系统给我们分配内存时，只需要分配两个页面大小就可以，分别存放代码和数据

调整段偏移

汇总所有符号

每个obj文件在编译时都会生成自己的符号表，我们要把这些符号都合并起来进行符号解析

完成符号的重定位

在进行合并段，调整段偏移时，输入文件的各个段在连接后的虚拟地址就已经确定了，这一步完成后，连接器开始计算各个符号的虚拟地址，因为各个符号在段内的相对位置是固定的，所以段内各个符号的地址也已经是确定的了，只不过连接器需要给每个符号加上一个偏移量，使他们能够调整到正确的虚拟地址，这就是符号的重定位过程

在 elf 文件中，有一个叫重定位表的结构专门用来保存这些与从定位有关的信息，重定位表在elf文件中往往是一个或多个段

5. 数据和指令

上面是代码编译链接的过程，得到了可执行的文件后，程序在内存中是如何运行的，数据和指令又分别是什么

所有的全局变量和静态变量都是数据，除此之外都是指令（包括**局部变量**）

虚拟地址空间

在每个程序运行的时候，我们的操作系统都会给他分配一个固定大小的虚拟地址空间（x86，32bit，Linux内核下默认大小为4G），那这段内存分配结构如下：

整个4G的空间有1G是供操作系统使用的内核空间，用户无法访问，还有3G是我们的用户空间，以供该虚拟地址空间上进程的运行，在这3G的用户空间中又被分成了很多段，从0地址开始的128M大小是系统的预留空间，用户也是无法访问的。

接下来是.text段，该段空间中存放的是代码，然后是.data段和.bss段，这两段里面存放的都是数据，但又有不同：data段中存放的数据是已经初始化并且初始化值不为0的数据，而.bss段中存放的是未经初始化或者初始化为0的数据（注：bss better save space(更好的节省空间)）

结语

经过以上分析，我们发现编译过程并不像想象的那么简单，而是要经过预处理、编译、汇编、链接。尽管我们平时使用gcc命令的时候没有关心中间结果，但每次程序的编译都少不了这几个步骤。也不用为上述繁琐过程而烦恼，因为你仍然可以：

```
1 $ gcc hello.c # 编译
2 $ ./a.out # 执行
```