

从入行开始就对编译器有着一种敬畏之心，虽然 c4 省去了链接及输出可执行文件过程，

但是通过对 c4 的学习还是能够大体了解到代码是如何被装载如内存再去执行的。

本文是对开源小项目 [c4](#) 进行学习的笔记，也可以视为对其实现的内存追踪，为的是能够更加深刻的了解程序实现的具体步骤。本项目的亮点在于仅仅通过 4 个函数就在内存中模拟出了对程序逐字逐句的解析并转换为伪汇编代码，最后在对伪汇编进行解析，输出程序执行结果；当然简洁性的代价就是该编译器仅支持特定语义的解析，并且刚才提到了编译和执行过程都是在内存中去模拟的，并不会生成特定阶段的目标文件比如：.i, .s, .out 文件，但是便于我们去阅读和学习，从中我们可以看出编译器最复杂的地方可能在于语义分析的逻辑部分了。

基本结构与函数的理解

- 全局变量

```
char *p, *lp, // 指向 source area 的指针，读入的代码会以文本信息的形式读
入在该区域
1 *data;      // data/bss 段指针
2 int *e, *le, // 指向 text area 的指针，解析代码后生成相应的伪汇编代码写
3 入该区域
4 *id,        // 指向当前正在设置或解析的语义关键字结构
5 *sym,       // 指向 symble area 的指针，能够解析的关键字、库函数以及函数
6 的参数和局部变量会存储在该区域
7 tk,        // 当前正在解析的代码的游标
8 ival,      // 当前解析的数字的值
9 ty,        // 语义关键字结构的类型
10 loc,      // 局部变量偏移量标识
11 line,     // 当前解析的代码行数
12 src,      // 出书解析的代码和生成的伪汇编代码的标识
debug;      // 输出汇编代码执行结果标识
```

- 枚举类型

```

1
2 // 类型、关键字、操作码标识
3 enum {
4     Num = 128, Fun, Sys, Glo, Loc, Id,
5     Char, Else, Enum, If, Int, Return, Sizeof, While,
6     Assign, Cond, Lor, Lan, Or, Xor, And, Eq, Ne, Lt, Gt, Le, Ge, Shl,
7     Shr, Add, Sub, Mul, Div, Mod, Inc, Dec, Brak
8 };
9
1 // 汇编操作码
0 enum
1 { LEA ,IMM ,JMP ,JSR ,BZ ,BNZ ,ENT ,ADJ ,LEV ,LI ,LC ,SI ,SC ,PSH
1
1 '
2
2 OR ,XOR ,AND ,EQ ,NE ,LT ,GT ,LE ,GE ,SHL ,SHR ,ADD ,SUB ,MUL ,D
1
1 IV ,MOD ,
3
3 OPEN,READ,CLOS,PRTF,MALC,MSET,MCMP,EXIT };
1
4
4 // 数据类型
1 enum { CHAR, INT, PTR };
5
1
1 // 语义关键字结构，方便进行偏移读取
6 enum { Tk, Hash, Name, Class, Type, Val, HClass, HType, HVal, Idsz };
1
7

```

- 函数

```

1/**
2 * 从 p 中读取下一个字符并进行处理
3 * 通过 tk 进行标识
4 * 遍历关键字结构体并用 id 指向该结构首地址
5 */
6void next();
7
8/**
9 * 操作符语义解析(如:sizeof, &, |, ++, --, =, ==等等)
10 * 生成伪汇编代码，写到 text area 中
11 * @param lev 操作符标识
12 */
13void expr(int lev);
14
15/**

```

```

16 * 逻辑关键字解析(如: if, while, return 等)
17 * 生成伪汇编代码, 写到 text area 中
18 */
19 void stmt();
20
21 /**
22 * 读取源代码文件
23 * 生成伪汇编代码
24 * 执行伪汇编代码
25 */
26 int main(int argc, char **argv)

```

准备阶段

为了便于调试我利用 `c4.c` 建立一个 Xcode 工程, 并且在 `main` 函数中给 `argc` 和 `argv` 进行了赋值:

```

int main(int argc, char **argv)
{
1  int fd, bt, ty, poolsz, *idmain;
2  int *pc, *sp, *bp, a, cycle; // vm registers
3  int i, *t; // temps
4
5  // 为输入参数赋值
6  argc = 3;
7  char **av = argv;
8  ++av;
9  *av = "-d";
10 ++av;
11 *av = "/Users/ronglei/Documents/Code/c4/add.c";
12
13 --argc; ++argv;
14 if (argc > 0 && **argv == '-' && (*argv)[1] == 's') { src = 1; --
15 argc; ++argv; }
16 if (argc > 0 && **argv == '-' && (*argv)[1] == 'd') { debug = 1; -
17 --argc; ++argv; }
18 if (argc < 1) { printf("usage: c4 [-s] [-d] file ...\n"); return -
19 1; }
20
...

```

在此我们以编译一下代码段为例子进行解析:

```

1#include <stdio.h>
2
3int add(int x, int y)
4{
5    int sum;
6    sum = x+y;
7    return sum;
8}
9
10void main()
11{
12    int a, b, c;
13    a = 1;
14    b = 2;
15    c = add(a, b);
16    printf("%d\n", c);
17}

```

接下来就是为各个区域申请内存空间：

```

    int main(int argc, char **argv)
    {
1      ...
2
3      // 申请内存空间
4      poolsz = 256*1024; // arbitrary size
5      if (!(sym = malloc(poolsz))) { printf("could not malloc(%d) symbol
6area\n", poolsz); return -1; }
7      if (!(le = e = malloc(poolsz))) { printf("could not malloc(%d)
8text area\n", poolsz); return -1; }
9      if (!(data = malloc(poolsz))) { printf("could not malloc(%d) data
10area\n", poolsz); return -1; }
11     if (!(sp = malloc(poolsz))) { printf("could not malloc(%d) stack
12area\n", poolsz); return -1; }
13
14     memset(sym, 0, poolsz);
15     memset(e, 0, poolsz);
16     memset(data, 0, poolsz);
17
18     if ((fd = open(*argv, 0)) < 0) { printf("could not open(%)s\n",
19*argv); return -1; }
20     if (!(lp = p = malloc(poolsz))) { printf("could not malloc(%d)
    source area\n", poolsz); return -1; }
        if ((i = (int)read(fd, p, poolsz-1)) <= 0) { printf("read()

```

```
returned %d\n", i); return -1; }
```

...

申请完成之后内存布局是这样的：



在申请完 symbol area, text area, data area, stack area 这四片区域之后，程序将支持的关键字和库函数转换为特定结构写入到 symbol 所对应的内存区域中，方便之后解析的时候使用，对应代码：

```
int main(int argc, char **argv)
1{
2    ...
3
4    p = "char else enum if int return sizeof while "
5    "open read close printf malloc memset memcmp exit void main";
6    i = Char; while (i <= While) { next(); id[Tk] = i++; } // add
7keywords to symbol table
8    i = OPEN; while (i <= EXIT) { next(); id[Class] = Sys; id[Type] =
9INT; id[Val] = i++; } // add library to symbol table
10    next(); id[Tk] = Char; // handle void type
11    next(); idmain = id; // keep track of main
12
    ...
}
```

执行完成后 symbol area 内存布局情况如图：

| | | | | | | | | | | | | |
|--------|------------|--------|------|------|-------|------|---------|--------|-------|------|--------|--|
| idmain | 0x0264/153 | Id | Hash | Name | Fun | CHAR | 0x40050 | 0 | 0 | 0 | main | |
| | 0x0240/144 | Char | Hash | Name | 0 | 0 | 0 | 0 | 0 | 0 | void | |
| | 0x021c/135 | Id | Hash | Name | Sys | INT | EXIT | 0 | 0 | 0 | exit | |
| | 0x01f8/126 | Id | Hash | Name | Sys | INT | MCMP | 0 | 0 | 0 | memcmp | |
| | 0x01d4/117 | Id | Hash | Name | Sys | INT | MSET | 0 | 0 | 0 | memset | |
| | 0x01b0/108 | Id | Hash | Name | Sys | INT | MALC | 0 | 0 | 0 | malloc | |
| | 0x018c/99 | Id | Hash | Name | Sys | INT | PRTF | 0 | 0 | 0 | printf | |
| | 0x0168/90 | Id | Hash | Name | Sys | INT | CLOS | 0 | 0 | 0 | close | |
| | 0x0144/81 | Id | Hash | Name | Sys | INT | READ | 0 | 0 | 0 | read | |
| | 0x0120/72 | Id | Hash | Name | Sys | INT | OPEN | 0 | 0 | 0 | open | |
| | 0x00fc/63 | While | Hash | Name | 0 | 0 | 0 | 0 | 0 | 0 | while | |
| | 0x00d8/54 | Sizeof | Hash | Name | 0 | 0 | 0 | 0 | 0 | 0 | sizeof | |
| | 0x00b4/45 | Return | Hash | Name | 0 | 0 | 0 | 0 | 0 | 0 | return | |
| | 0x0090/36 | Int | Hahs | Name | 0 | 0 | 0 | 0 | 0 | 0 | int | |
| | 0x006c/27 | If | Hash | Name | 0 | 0 | 0 | 0 | 0 | 0 | if | |
| | 0x0048/18 | Enum | Hash | Name | 0 | 0 | 0 | 0 | 0 | 0 | enum | |
| | 0x0024/9 | Else | Hash | Name | 0 | 0 | 0 | 0 | 0 | 0 | else | |
| | 0x0000/0 | Char | Hash | Name | 0 | 0 | 0 | 0 | 0 | 0 | char | |
| sym | 0x0000/0 | Tk | Hash | Name | Class | Type | Val | HClass | HType | HVal | | |

其中 Hash 段和 Name 段是通过关键字或函数名称字符串运算得到到，拿关键字 char 举例，源码这里存在一些疏漏，我们会在最后说明：

1Hash 字段值为: `(c+h+a+r)*147*64+strlen(char)`
 2Name 字段值为: `"*(int *)char..."`

语义解析阶段

接下来就要进行语义的解析并在 text area 生成响应的伪汇编代码：

```

1int main(int argc, char **argv)
2{
3    ...
4
5    // parse declarations
6    line = 1;
7    next();
8    // 检测到 int add(...)
9    while (tk) {
10        bt = INT; // basetype
11        if (tk == Int) next(); // 接下来处理 add 函数 id 指向 add 所构建的
12结构体 tk=Id
13        else if (tk == Char) { next(); bt = CHAR; }
14        else if (tk == Enum) {
15            next();
16            if (tk != '{') next();

```

```

17         if (tk == '{') {
18             next();
19             i = 0;
20             while (tk != '}') {
21                 if (tk != Id) { printf("%d: bad enum
22identifier %d\n", line, tk); return -1; }
23                 next();
24                 if (tk == Assign) {
25                     next();
26                     if (tk != Num) { printf("%d: bad enum
27initializer\n", line); return -1; }
28                     i = ival;
29                     next();
30                 }
31                 id[Class] = Num; id[Type] = INT; id[Val] = i++;
32                 if (tk == ',') next();
33             }
34             next();
35         }
36     }
37     while (tk != ';' && tk != '}') {
38         ty = bt;
39         while (tk == Mul) { next(); ty = ty + PTR; }
40         if (tk != Id) { printf("%d: bad global declaration\n",
41line); return -1; }
42         if (id[Class]) { printf("%d: duplicate global
43definition\n", line); return -1; }
44         next(); // 过滤掉 '('
45         id[Type] = ty;
46         if (tk == '(') { // function
47             id[Class] = Fun;
48//             id[Val] = (int)(e + 1);
49             id[Val] = (long)(e + 1);
50             next(); i = 0;
51             while (tk != ')') { // 处理 add 函数的参数列表
52                 ty = INT;
53                 if (tk == Int) next();
54                 else if (tk == Char) { next(); ty = CHAR; }
55                 while (tk == Mul) { next(); ty = ty + PTR; }
56                 if (tk != Id) { printf("%d: bad parameter
57declaration\n", line); return -1; }
58                 if (id[Class] == Loc) { printf("%d: duplicate
59parameter definition\n", line); return -1; }
60                 id[HClass] = id[Class]; id[Class] = Loc;

```

```

61         id[HType] = id[Type]; id[Type] = ty;
62         id[HVal] = id[Val]; id[Val] = i++;
63         next();
64         if (tk == ',') next();
65     }
66     next(); // add 函数参数解析完成
67     if (tk != '{') { printf("%d: bad function
68definition\n", line); return -1; }
69     loc = ++i;
70     next();
71     while (tk == Int || tk == Char) { // 处理函数内的局部变
72量
73         bt = (tk == Int) ? INT : CHAR;
74         next();
75         while (tk != ';') {
76             ty = bt;
77             while (tk == Mul) { next(); ty = ty + PTR; }
78             if (tk != Id) { printf("%d: bad local
79declaration\n", line); return -1; }
80             if (id[Class] == Loc) { printf("%d: duplicate
81local definition\n", line); return -1; }
82             id[HClass] = id[Class]; id[Class] = Loc;
83             id[HType] = id[Type]; id[Type] = ty;
84             id[HVal] = id[Val]; id[Val] = ++i;
85             next();
86             if (tk == ',') next();
87         }
88         next();
89     }
90     *++e = ENT; *++e = i - loc; // 文本段 ENT 入口 (i-loc)局
91部变量个数
92     while (tk != '}') stmt(); // 处理函数内部的操作
93     *++e = LEV;
94     id = sym; // 将函数内部的局部变量的 Class Type Val 清空
95     while (id[Tk]) {
96         if (id[Class] == Loc) {
97             id[Class] = id[HClass];
98             id[Type] = id[HType];
99             id[Val] = id[HVal];
100         }
101         id = id + Idsz;
102     }
    }
    else {

```



```
        id[Class] = Glo;
        id[Val] = (int)data;
        data = data + sizeof(int);
    }
    if (tk == ',') next();
}
next();
}

...

```

再此过程中, 会为曾将解析过的函数、函数的参数以及函数作用于内声明的变量, 在 symbol area 内生成相应的结构来对他们进行标识和记录, 最后在对它们的特定字段进行清理, 如图:

| | | | | | | | | | | | | |
|-------------------|---|--------|------|------|--------|--------|---------|--------|-------|------|--------|------|
| 0x0360/216 | → | Id | Hash | Name | Loc->0 | INT->0 | 4->0 | 0 | 0 | 0 | c | 局部变量 |
| 0x033c/207 | → | Id | Hash | Name | Loc->0 | INT->0 | 3->0 | 0 | 0 | 0 | b | |
| 0x0318/198 | → | Id | Hash | Name | Loc->0 | INT->0 | 2->0 | 0 | 0 | 0 | a | |
| 0x02f4/189 | → | Id | Hash | Name | Loc->0 | INT->0 | 4->0 | 0 | 0 | 0 | sum | |
| 0x02d0/180 | → | Id | Hash | Name | Loc->0 | INT->0 | 1->0 | 0 | 0 | 0 | y | |
| 0x02ac/171 | → | Id | Hash | Name | Loc->0 | INT->0 | 0->0 | 0 | 0 | 0 | x | |
| 0x0288/162 | → | Id | Hash | Name | Fun | INT | 0x40004 | 0 | 0 | 0 | add | |
| idmain 0x0264/153 | → | Id | Hash | Name | Fun | CHAR | 0x40050 | 0 | 0 | 0 | main | 库函数 |
| 0x0240/144 | → | Char | Hash | Name | 0 | 0 | 0 | 0 | 0 | 0 | void | |
| 0x021c/135 | → | Id | Hash | Name | Sys | INT | EXIT | 0 | 0 | 0 | exit | |
| 0x01f8/126 | → | Id | Hash | Name | Sys | INT | MCMP | 0 | 0 | 0 | memcmp | |
| 0x01d4/117 | → | Id | Hash | Name | Sys | INT | MSET | 0 | 0 | 0 | memset | |
| 0x01b0/108 | → | Id | Hash | Name | Sys | INT | MALC | 0 | 0 | 0 | malloc | |
| 0x018c/99 | → | Id | Hash | Name | Sys | INT | PRTF | 0 | 0 | 0 | printf | |
| 0x0168/90 | → | Id | Hash | Name | Sys | INT | CLOS | 0 | 0 | 0 | close | 关键字 |
| 0x0144/81 | → | Id | Hash | Name | Sys | INT | READ | 0 | 0 | 0 | read | |
| 0x0120/72 | → | Id | Hash | Name | Sys | INT | OPEN | 0 | 0 | 0 | open | |
| 0x00fc/63 | → | While | Hash | Name | 0 | 0 | 0 | 0 | 0 | 0 | while | |
| 0x00d8/54 | → | Sizeof | Hash | Name | 0 | 0 | 0 | 0 | 0 | 0 | sizeof | |
| 0x00b4/45 | → | Return | Hash | Name | 0 | 0 | 0 | 0 | 0 | 0 | return | |
| 0x0090/36 | → | Int | Hahs | Name | 0 | 0 | 0 | 0 | 0 | 0 | int | |
| 0x006c/27 | → | If | Hash | Name | 0 | 0 | 0 | 0 | 0 | 0 | if | 关键字 |
| 0x0048/18 | → | Enum | Hash | Name | 0 | 0 | 0 | 0 | 0 | 0 | enum | |
| 0x0024/9 | → | Else | Hash | Name | 0 | 0 | 0 | 0 | 0 | 0 | else | |
| sym 0x0000/0 | → | Char | Hash | Name | 0 | 0 | 0 | 0 | 0 | 0 | char | |
| | | Tk | Hash | Name | Class | Type | Val | HClass | HType | HVal | | |

在 text area 生成的伪汇编代码后的内存布局情况如图:

| | | | | | | | | | | |
|------------|---|---------|-----|------|---------|-----|---------|---------|-----|-----|
| | | LI | PHS | PRTF | ADJ | 2 | LEV | 0 | 0 | 0 |
| 0x400d8 | → | 0x40004 | ADJ | 2 | SI | IMM | %d\n | PSH | LEA | -3 |
| 0x400b4 | → | LI->PSH | LEA | -1 | LI | PSH | LEA | -2 | LI | JSR |
| 0x40090 | → | SI | LEA | -2 | LI->PSH | IMM | 2 | SI | LEA | -3 |
| 0x4006c | → | LEV | LEV | ENT | 3 | LEA | -1 | LI->PSH | IMM | 1 |
| 0x40048 | → | PSH | LEA | 2 | LI | ADD | SI | LEA | -1 | LI |
| 0x40024 | → | 0 | ENT | 1 | LEA | -1 | LI->PSH | LEA | 3 | LI |
| le 0x40000 | → | | | | | | | | | |

汇编执行阶段

最后就是对 text area 生成的汇编代码进行解析和执行，代码如下：

```

1 int main(int argc, char **argv)
2 {
3     ...
4
5     // 获取 main 函数入口地址赋值给 pc
6     if (!(pc = *(long *) (idmain+Val))) { printf("main() not
7 defined\n"); return -1; }
8
9     if (src) return 0;
1    // 设置栈区域，栈是从高位地址向低位地址扩展的
0    sp = (int *) ((char *) sp + poolsz);
1    *--sp = EXIT; // call exit if main returns
1    *--sp = PSH; t = sp;
1    *--sp = argc;
2    *--sp = (int) argv;
1    *--sp = (int) t;
3
1    // run...
4    cycle = 0;
1    while (1) {
5        i = *pc++; ++cycle;
1        if (debug) {
6            printf("%d> %.4s", cycle,
1
7 &"LEA ,IMM ,JMP ,JSR ,BZ ,BNZ ,ENT ,ADJ ,LEV ,LI ,LC ,SI ,SC ,PSH
1, "
8
1 "OR ,XOR ,AND ,EQ ,NE ,LT ,GT ,LE ,GE ,SHL ,SHR ,ADD ,SUB ,MUL ,
9 DIV ,MOD , "
2
10 "OPEN,READ,CLOS,PRTF,MALC,MSET,MCMP,EXIT,"[i * 5]);

```

```

0         if (i <= ADJ) printf(" %d\n", *pc); else printf("\n");
2     }
1     if      (i == LEA) a = (int)(bp + *pc++);
2// load local address
2         else if (i == IMM) a = *pc++;
2// load global address or immediate
3         else if (i == JMP) pc = (int *)*pc;
2// jump
4         else if (i == JSR) { *--sp = (int)(pc + 1); pc = (int *)*pc; }
2// jump to subroutine
5         else if (i == BZ)  pc = a ? pc + 1 : (int *)*pc;
2// branch if zero
6         else if (i == BNZ) pc = a ? (int *)*pc : pc + 1;
2// branch if not zero
7         else if (i == ENT) { *--sp = (int)bp; bp = sp; sp = sp -
2*pc++; } // enter subroutine
8         else if (i == ADJ) sp = sp + *pc++;
2// stack adjust
9         else if (i == LEV) { sp = bp; bp = (int *)*sp++; pc = (int
3*)*sp++; } // leave subroutine
0         else if (i == LI)  a = *(int *)a;
3// load int
1         else if (i == LC)  a = *(char *)a;
3// load char
2         else if (i == SI)  *(int *)*sp++ = a;
3// store int
3         else if (i == SC)  a = *(char *)*sp++ = a;
3// store char
4         else if (i == PSH) *--sp = a;
3// push
5
3         else if (i == OR)  a = *sp++ |  a;
6         else if (i == XOR) a = *sp++ ^  a;
3         else if (i == AND) a = *sp++ &  a;
7         else if (i == EQ)  a = *sp++ == a;
3         else if (i == NE)  a = *sp++ != a;
8         else if (i == LT)  a = *sp++ <  a;
3         else if (i == GT)  a = *sp++ >  a;
9         else if (i == LE)  a = *sp++ <= a;
4         else if (i == GE)  a = *sp++ >= a;
0         else if (i == SHL) a = *sp++ << a;
4         else if (i == SHR) a = *sp++ >> a;
1         else if (i == ADD) a = *sp++ +  a;
4         else if (i == SUB) a = *sp++ -  a;

```

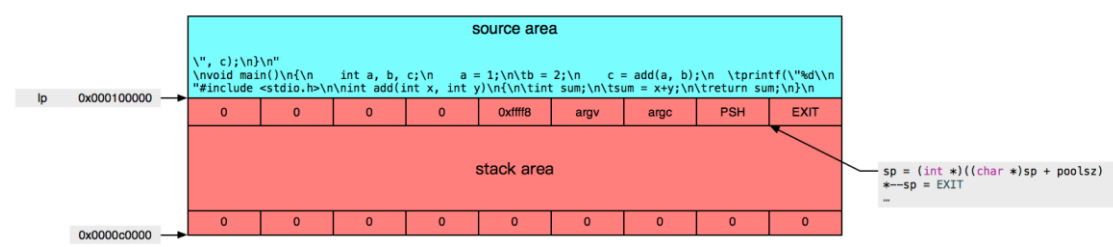
```

2     else if (i == MUL) a = *sp++ * a;
4     else if (i == DIV) a = *sp++ / a;
3     else if (i == MOD) a = *sp++ % a;
4
4     else if (i == OPEN) a = open((char *)sp[1], *sp);
4     else if (i == READ) a = read(sp[2], (char *)sp[1], *sp);
5     else if (i == CLOS) a = close(*sp);
4     else if (i == PRTF) { t = sp + pc[1]; a = printf((char *)t[-
61], t[-2], t[-3], t[-4], t[-5], t[-6]); }
4     else if (i == MALC) a = (int)malloc(*sp);
7     else if (i == MSET) a = (int)memset((char *)sp[2], sp[1],
4*sp);
8     else if (i == MCMP) a = memcmp((char *)sp[2], (char *)sp[1],
4*sp);
9     else if (i == EXIT) { printf("exit(%d) cycle = %d\n", *sp,
5cycle); return *sp; }
0     else { printf("unknown instruction = %d! cycle = %d\n", i,
5cycle); return -1; }
1 }
5}
2
5
3
5
4
5
5
5
5
6
5
7
5
8
5
9
6
0
6
1
6
2
6
3
6

```

| |
|---|
| 4 |
| 6 |
| 5 |
| 6 |
| 6 |
| 6 |
| 7 |
| 6 |
| 8 |
| 6 |
| 9 |
| 7 |
| 0 |

这个阶段利用 `pc` 读取 `text area` 生成的伪汇编代码，然后利用指向 `stack area` 的指针 `sp` 模拟执行，示意图如下：



纠错

- 查找支持的关键字或库函数过程中，在 `next()` 函数以下代码段中：

```
1void next()
2{
3...
4  else if ((tk >= 'a' && tk <= 'z') || (tk >= 'A' && tk <=
5'Z') || tk == '_') {
6    pp = p - 1;
7    while ((*p >= 'a' && *p <= 'z') || (*p >= 'A' && *p <= 'Z')
8|| (*p >= '0' && *p <= '9') || *p == '_')
9      tk = tk * 147 + *p++;
10   tk = (tk << 6) + (p - pp);
11   id = sym;
12   while (id[Tk]) {
13     if (tk == id[Hash] && !memcmp((char *)id[Name], pp, p -
```

```

14pp)) { tk = id[Tk]; return; }
15     id = id + Idsz;
16 }
17 id[Name] = (int)pp;
18 id[Hash] = tk;
19 tk = id[Tk] = Id;
20 return;
    }
    ...
}

```

1 从代码中我们可以看出，在对之前添加的关键字或函数进行匹配时，是通过对
2 Hash 段和 Name 段进行比较得出的，条件表达式为即：

// 这条语句在编译的时候编译器会报错

```
tk == id[Hash] && !memcmp((char *)id[Name], pp, p - pp)
```

// 正确的写法

```
tk == id[Hash] && !memcmp((char *)id+Name), pp, p - pp)
```

1 当我们将 Name 进行存储的时候是将字符串指针强转为 int 型去存储的，即：

```
id[Name] = (int)pp;
```

1 这里面作者想表达的意思是，将字符串指针强转为整形指针存入到 Name 字段中，所以对
2 Name 字段进行存储是应该做如下修改：

```
id[Name] = *(int *)pp;
```

1 注意由于 char 类型占 1 个字节，int 类型占 4 个字节，所以此时 Name 段记录的字符
2 为 pp 所指向字符的前 4 个字节：
3 如果 pp = "char else..."，则 Name 字段内容转换为字符串为 "char"；
4 如果 pp = "if else"，则 Name 字段内容转换为字符串为 "if e"；
5 如果 pp = "return else"，则 Name 字段内容转换为字符串为 "retu"；
6 因此在该字段内容进行比较时，对比较的长度的要求就比较特殊，上面条件判断就应该
改为：

// 如果比较的字符串长度小于 4 个字节，比较长度就取该长度；

// 如果比较的字符串长度大于 4 个字节，比较长度就取 4 个字节长度进行比较，

因为在存储的之后我们只记录了字符串的前 4 个字节；

```
int cmp_len = (p - pp) > 4 ? 4 : (p - pp);

if (tk == id[Hash] && !memcmp((char *) (id + Name), pp, cmp_len))

{

    tk = id[Tk];

    return;

}
```

1 * 在对生成的汇编代码进行模拟运行的时候，需要获取之前存储的 `main` 函数以及自定义
2 函数的地址进行跳转，而当时存储的时候代码如下：

```
id[Val] = (int)(e + 1);
```

1 我们看到是将地址值强转成 `int` 型去存储的，对于 32 位机器来说，地址是通过 8 位 16
2 进制数去表示的，地址的长度正好是 `int` 型的长度 4 个字节；但在 64 位机器上，地址是
通过 16 位 16 进制数表示，地址长度为 8 个字节，此处如果想用 `int` 型去存储地址的
话，高 4 位地址会被截取，存储的只是低 4 位地址的内容，在跳转的时候就会出问题。