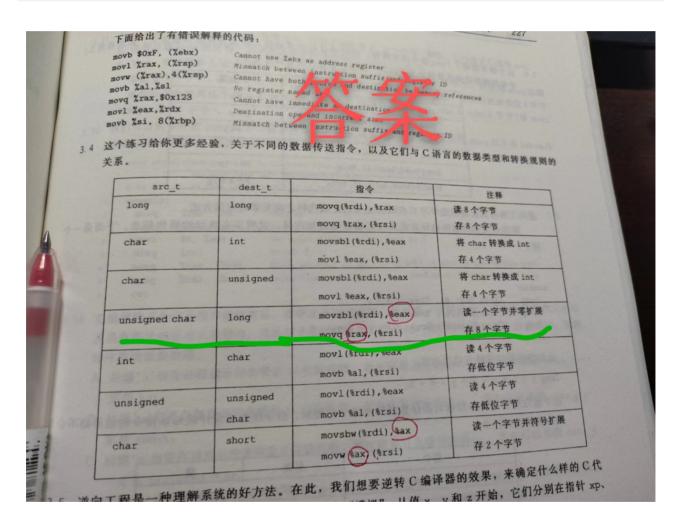
源码与汇编代码交织

有位同学问了一个很好的问题,我来试图重复下这个问题。

Q: 如图所示,从 unsigned char 通过数据类型转换到 long 为什么要这么写,我能写成下面的形式吗?

```
1 movzbq (%rdi), %rax
2 movq %rax, (%rsi)
```



当我看到这个问题的时候, 我觉得问的很棒, 于是我开始思考(时间滴答滴答流逝了...

我当时的回答是: 我觉得彳亍,但是我又没有证据说为什么行,于是我通过以下方法试图还原作者当时的想法。我写了下面一段 c 程序,并将其保存为 test.c

```
1 #include <stdio.h>
2 int main()
3 {
4    unsigned char c = 'c';
5    long l = (long)c;
6    return 0;
7 }
```

很明显,这是一段 c 代码,我想是否有办法能够看到这段代码的汇编呢? 我想到了 objdump 工具。 干是我做了下面的操作

```
1 man objdump
```

我查到了通过 -S 参数可以使得交织汇编和源代码, 欣喜若狂地执行下面指令

```
1 gcc -00 -g -o test test.c
2 objdump -S test > test.asm
```

接着使用 vim 打开 test.asm 文件, 我看到了下面的内容:

```
1 #include <stdio.h>
2 int main()
3 {
4 1129: f3 Of 1e fa endbr64
    112d: 55
5
                             push %rbp
    112e: 48 89 e5
                              mov %rsp,%rbp
6
    unsigned char c = 'c';
7
    1131: c6 45 f7 63
                              movb $0x63,-0x9(%rbp)
8
    long l = (long)c;
    1135: 0f b6 45 f7
                              movzbl -0x9(%rbp),%eax
10
11
    1139: 48 89 45 f8
                              mov %rax,-0x8(%rbp)
12
    return 0;
    113d: b8 00 00 00 00
                              mov $0x0,%eax
13
14 }
```

哦吼,这就是我想要的,我们看到 long l = (long)c; 下面的两句汇编实际上就是执行了和问题图中一样的操作,很很明显 eax 是 32 位寄存器。上面的实验是在 amd64 上做的,刚好我还有一台 arm64 的 ubuntu 机器,做了下面操作:

```
parallels@jpzhu-db-edu:~$ cat /proc/version
Linux version 5.4.0-126-generic (build@bos02-arm64-060) (gcc version 9.4.0 (Ubuntu 9.4.0-1ubuntu1~20.04.1)) #142-Ubuntu SMP Fri Aug 26 12:15:55 UTC 2022
parallels@jpzhu-db-edu:~$ cat test.c
#include <stdio.h>

int main()
{
    unsigned char c = 'c';
    long l = (long)c;
    return 0;
}
parallels@jpzhu-db-edu:~$ gcc -00 -g -o test test.c
parallels@jpzhu-db-edu:~$ objdump -S test > test.asm
parallels@jpzhu-db-edu:~$ vim test.asm
```

我得到了如下输出, arm64 下的 ldrb 指令也是扩展为 32 位后,在 store 回去,gcc 编译出的结果看上去都是经历了 unsigned char -> int -> long 的过程。

```
1 #include <stdio.h>
2 int main()
3 {
  71c: d10043ff sub
                                sp, sp, \#0 \times 10
      unsigned char c = 'c';
6 720: 52800c60
                         mov
                                 w0, #0x63
                                                               // #99
7 724: 39001fe0
                                 w0, [sp, #7]
                         strb
      long l = (long)c;
  728: 39401fe0
                         ldrb
                                 w0, [sp, #7]
10 72c: f90007e0
                                 x0, [sp, #8]
                        str
     return 0;
11
12 730: 52800000
                                                                // #0
                                 w0, \#0\times0
                         mov
13 }
```

于是我通过真实的 gcc 编译出的汇编代码,接受了如图所示的这种写法。

那么还有一个问题,为什么 gcc 在编译代码的时候要选择这种方法呢?这也是当时这位提出问题的学生追问我的。我再次陷入了思考......

我检索到了 https://en.wikipedia.org/wiki/Sign_extension ,我发现在 zero extension 部分,通常情况下,在 64 位的机器上,都是先扩展到 32 位,同时高位的 32 会也会被置零,这是一个传统。我想这个传统恐怕是为了最大程度和 32 位机器兼容.