

/dev/mem

/dev/kmem

/dev/null

/dev/zero

/dev/port

/dev/full

/dev/random

/dev/urandom

/dev/kmsg

/dev/mem --- physical memory, read / write

/dev/kmem --- the kernel virtual memory rather than physical memory is accessed

/dev/full --- Writes to the /dev/full device will fail with an ENOSPC error.

This can be used to test how a program handles disk-full errors.

Reads from the /dev/full device will return \0 characters.

Seeks on /dev/full will always succeed.

/dev/port --- port is similar to mem, but the I/O ports are accessed.

The difference between /dev/random and /dev/urandom

When the entropy pool is empty, reads from `/dev/random` will block until additional environmental noise is gathered.

A read from the `/dev/urandom` device will not block waiting for more entropy.

in `drivers/char/mem.c`

```

1. static const struct memdev {
2.     const char *name;
3.     umode_t mode;
4.     const struct file_operations *fops;
5.     struct backing_dev_info *dev_info;
6. } devlist[] = {
7.     [1] = { "mem", 0, &mem_fops, &directly_mappable_cdev_bdi },
8. #ifdef CONFIG_DEVMEM
9.     [2] = { "kmem", 0, &kmem_fops, &directly_mappable_cdev_bdi },
10. #endif
11.     [3] = { "null", 0666, &null_fops, NULL },
12. #ifdef CONFIG_DEVPORT
13.     [4] = { "port", 0, &port_fops, NULL },
14. #endif
15.     [5] = { "zero", 0666, &zero_fops, &zero_bdi },
16.     [7] = { "full", 0666, &full_fops, NULL },
17.     [8] = { "random", 0666, &random_fops, NULL },
18.     [9] = { "urandom", 0666, &urandom_fops, NULL },
19. #ifdef CONFIG_PRINTK
20.     [11] = { "kmsg", 0644, &kmsg_fops, NULL },
21. #endif
22. };
23.
24. ....
25.
26. static int __init chr_dev_init(void)
27. {
28.     int minor;
29.     int err;
30.
31.     err = bdi_init(&zero_bdi);
32.     if (err)
33.         return err;
34.
35.     if (register_chrdev(MEM_MAJOR, "mem", &memory_fops))
36.         printk("unable to get major %d for memory devs\n", MEM_MAJOR);
37.
38.     mem_class = class_create(THIS_MODULE, "mem");
39.     if (IS_ERR(mem_class))
40.         return PTR_ERR(mem_class);
41.
42.     mem_class->devnode = mem_devnode;
43.     for (minor = 1; minor < ARRAY_SIZE(devlist); minor++) {
44.         if (!devlist[minor].name)
45.             continue;
46.
47.         /*
48.          * Create /dev/port?
49.          */
50.         if ((minor == DEVPORT_MINOR) && !arch_has_dev_port())
51.             continue;
52.
53.         device_create(mem_class, NULL, MKDEV(MEM_MAJOR, minor),

```

```
54.             NULL, devlist[minor].name);
55.     }
56.
57.     return tty_init();
58. }
```

/dev/mem

"/dev/mem" 对应的file_operations is

```
1.  static const struct file_operations mem_fops = {
2.      .llseek      = memory_llseek,
3.      .read         = read_mem,
4.      .write        = write_mem,
5.      .mmap         = mmap_mem,
6.      .open         = open_mem,
7.      .get_unmapped_area = get_unmapped_area_mem,
8.  };
```

```

1.  /*
2.  * This function reads the *physical* memory. The f_pos points directly to the
3.  * memory location.
4.  */
5.  static ssize_t read_mem(struct file *file, char __user *buf,
6.                          size_t count, loff_t *ppos)
7.  {
8.      phys_addr_t p = *ppos;
9.      ssize_t read, sz;
10.     char *ptr;
11.
12.     if (p != *ppos)
13.         return 0;
14.
15.     if (!valid_phys_addr_range(p, count)) ①
16.         return -EFAULT;
17.     read = 0;
18. #ifdef __ARCH_HAS_NO_PAGE_ZERO_MAPPED
19.     /* we don't have page 0 mapped on sparc and m68k.. */
20.     if (p < PAGE_SIZE) {
21.         sz = size_inside_page(p, count);
22.         if (sz > 0) {
23.             if (clear_user(buf, sz))
24.                 return -EFAULT;
25.             buf += sz;
26.             p += sz;
27.             count -= sz;
28.             read += sz;
29.         }
30.     }
31. #endif
32.
33.     while (count > 0) {
34.         unsigned long remaining;
35.
36.         sz = size_inside_page(p, count);
37.
38.         if (!range_is_allowed(p >> PAGE_SHIFT, count))
39.             return -EPERM;
40.
41.         /*
42.          * On ia64 if a page has been mapped somewhere as uncached, then
43.          * it must also be accessed uncached by the kernel or data
44.          * corruption may occur.
45.          */
46.         ptr = xlate_dev_mem_ptr(p); ②
47.         if (!ptr)
48.             return -EFAULT;
49.
50.         remaining = copy_to_user(buf, ptr, sz); ③
51.         unxlate_dev_mem_ptr(p, ptr);
52.         if (remaining)
53.             return -EFAULT;

```

```

54.
55.         buf += sz;
56.         p += sz;
57.         count -= sz;
58.         read += sz;
59.     }
60.
61.     *ppos += read;
62.     return read;
63. }

```

①

check 读的区域是否在physical memory之内

```

1.  static inline int valid_phys_addr_range(phys_addr_t addr, size_t count)
2.  {
3.      return addr + count <= __pa(high_memory);
4.  }

```

high_memory是memory的最高边界的virtual address

②

in arch/arm/include/asm/io.h

```

1.  /*
2.   * Convert a physical pointer to a virtual kernel pointer for /dev/mem
3.   * access
4.   */
5.  #define xlate_dev_mem_ptr(p)    __va(p)

```

p是/dev/mem file offset ==> physical address

转变成virtual address

③

ptr is virtual address

/dev/kmem

"/dev/kmem" 对应的file_operations is

```
1. static const struct file_operations kmem_fops = {  
2.     .llseek      = memory_llseek,  
3.     .read         = read_kmem,  
4.     .write        = write_kmem,  
5.     .mmap         = mmap_kmem,  
6.     .open         = open_kmem,  
7.     .get_unmapped_area = get_unmapped_area_mem,  
8. };
```

/dev/kmem对应的文件是kernel看到的address space。大致如下

vmalloc : 0xf0000000 - 0xff000000 (240 MB)

lowmem : 0xc0000000 - 0xef800000 (760 MB)

pkmap : 0xbfe00000 - 0xc0000000 (2 MB)

modules : 0xbf000000 - 0xbfe00000 (14 MB)

.text : 0xc0008000 - 0xc05a6d14 (5756 kB)

.init : 0xc05a7000 - 0xc05da000 (204 kB)

.data : 0xc05da000 - 0xc06098c0 (191 kB)

.bss : 0xc06098c0 - 0xc0679c74 (449 kB)

即kernel看到的space是从0xc000,0000(virtual address , 对应physical address 0)开始 ,
到high_memory , 是对应到physical memory的kernel space , 在此上还有vmalloc space , 这
并不象lowmem那样直接对应到physical memory,它的physical memory也来自与0xc000,0000到
high_memory。


```

1.  /*
2.  * This function reads the *virtual* memory as seen by the kernel.
3.  */
4.  static ssize_t read_kmem(struct file *file, char __user *buf,
5.                          size_t count, loff_t *ppos)
6.  {
7.      unsigned long p = *ppos;                ❶
8.      ssize_t low_count, read, sz;
9.      char *kbuf; /* k-addr because vread() takes vmlist_lock rwlock */
10.     int err = 0;
11.
12.     read = 0;
13.     if (p < (unsigned long) high_memory) {    ❷
14.         low_count = count;
15.         if (count > (unsigned long) high_memory - p)
16.             low_count = (unsigned long) high_memory - p;
17.
18. #ifdef __ARCH_HAS_NO_PAGE_ZERO_MAPPED
19.     /* we don't have page 0 mapped on sparc and m68k.. */
20.     if (p < PAGE_SIZE && low_count > 0) {
21.         sz = size_inside_page(p, low_count);
22.         if (clear_user(buf, sz))
23.             return -EFAULT;
24.         buf += sz;
25.         p += sz;
26.         read += sz;
27.         low_count -= sz;
28.         count -= sz;
29.     }
30. #endif
31.     while (low_count > 0) {
32.         sz = size_inside_page(p, low_count);
33.
34.         /*
35.          * On ia64 if a page has been mapped somewhere as
36.          * uncached, then it must also be accessed uncached
37.          * by the kernel or data corruption may occur
38.          */
39.         kbuf = xlate_dev_kmem_ptr((char *)p);    ❸
40.
41.         if (copy_to_user(buf, kbuf, sz))          ❹
42.             return -EFAULT;
43.         buf += sz;
44.         p += sz;
45.         read += sz;
46.         low_count -= sz;
47.         count -= sz;
48.     }
49. }
50.
51. if (count > 0) {                                ❺
52.     kbuf = (char *)__get_free_page(GFP_KERNEL);  ❻
53.     if (!kbuf)

```

```

54.         return -ENOMEM;
55.     while (count > 0) {
56.         sz = size_inside_page(p, count);
57.         if (!is_vmalloc_or_module_addr((void *)p)) { ⑦
58.             err = -ENXIO;
59.             break;
60.         }
61.         sz = vread(kbuf, (char *)p, sz); ⑧
62.         if (!sz)
63.             break;
64.         if (copy_to_user(buf, kbuf, sz)) { ⑨
65.             err = -EFAULT;
66.             break;
67.         }
68.         count -= sz;
69.         buf += sz;
70.         read += sz;
71.         p += sz;
72.     }
73.     free_page((unsigned long)kbuf);
74. }
75. *ppos = p;
76. return read ? read : err;
77. }

```

①

p是/dev/kmem file的offset，它的size就可以超过physical memory size，可以达vmalloc space的边界

②

表示要读取的文件区域的上边界是小于high_memory的，即落在physical memory之内

③

in arch/arm/include/asm/io.h

```
#define xlate_dev_kmem_ptr(p) p
```

由于这里p即是kernel address , 所以无需转换

④

copy小于high_memory的data

⑤

count > 0

表示读取区域的下边界大于high_memory , 即超出了0xc000,0000 + physical memory , 也就是
vmalloc区域

⑥

读取vmalloc区域 , 需要allocate one page 过度一下 (Why???)

⑦

check p的合法性 , 超过vmalloc区域也是非法的

⑧

读取到过度page中

⑨

复制到用户空间

/dev/null

"/dev/null" 对应的file_operations is

```
1. static const struct file_operations null_fops = {
2.     .llseek      = null_llseek,
3.     .read        = read_null,
4.     .write       = write_null,
5.     .aio_read    = aio_read_null,
6.     .aio_write   = aio_write_null,
7.     .splice_write = splice_write_null,
8. };
9.
10.
11. static ssize_t read_null(struct file *file, char __user *buf,
12.                          size_t count, loff_t *ppos)
13. {
14.     return 0;
15. }
16.
17. static ssize_t write_null(struct file *file, const char __user *buf,
18.                           size_t count, loff_t *ppos)
19. {
20.     return count;
21. }
```

/dev/full

"/dev/full" 对应的file_operations is

```
1. static const struct file_operations full_fops = {
2.     .llseek      = full_llseek,
3.     .read        = new_sync_read,
4.     .read_iter   = read_iter_zero,
5.     .write       = write_full,
6. };
7.
8. static ssize_t write_full(struct file *file, const char __user *buf,
9.                           size_t count, loff_t *ppos)
10. {
11.     return -ENOSPC;
12. }
```

/dev/port

"/dev/port" 对应的file_operations is

```
1. static const struct file_operations port_fops = {  
2.     .llseek      = memory_llseek,  
3.     .read        = read_port,  
4.     .write       = write_port,  
5.     .open        = open_port,  
6. };
```

???