

Paging On Linux x86-64

Jay Khandkar

26 January 2021

Abstract

In this document we shall take a look at how exactly Linux handles paging by writing a character driver that accepts virtual addresses by means of the `ioctl` interface and walks through the page tables of the process to retrieve the corresponding physical address. This hands-on approach shall demonstrate how to read and understand the kernel source, as well as highlight some important aspects of paging, such as the distinction between page frame numbers and page frames, pages and page frames etc. All aspects of paging described in this document will be specific to the **x86-64** platform.

1 The Virtual Address Space

Firstly, it may be noted that starting with kernel version 4.11, Linux uses five-level page tables as opposed to four-level before. However, for the sake of simplicity, we shall work with four-level tables when describing aspects of paging, and the concept can easily be extended to five-level tables, since it is just one more level of indirection.

Let us now look at how the virtual address space is split in Linux. All 64 bits of the virtual address space are not used. For four-level tables, only the first 48 are used, and bits 48:63 are either all zero (in user space) or all one (in kernel space). To see in detail how this is implemented, let us consult

`Documentation/x86/x86_64/mm.rst`:

Start addr	Offset	End addr	Size	VM area description
0000000000000000	0	00007fffffffffff	128 TB	user-space virtual memory, different per mm
0000800000000000	+128 TB	ffff7fffffffffff	~16M TB	... huge, almost 64 bits wide hole of non-canonical virtual memory addresses up to the -128 TB starting offset of kernel mappings.
Kernel-space virtual memory, shared between all processes:				
ffff800000000000	-128 TB	ffff87ffffffffff	8 TB	... guard hole, also reserved for hypervisor
ffff800000000000	-120 TB	ffff887fffffffff	0.5 TB	LDT remap for PTI
ffff880000000000	-119.5 TB	ffffc87fffffffff	64 TB	direct mapping of all physical memory (page_offset_base)
ffffc80000000000	-55.5 TB	ffffc87fffffffff	0.5 TB	... unused hole
ffffc90000000000	-55 TB	ffffc87fffffffff	32 TB	vmalloc/ioremap space (vmalloc_base)
ffffc90000000000	-23 TB	ffffc97fffffffff	1 TB	... unused hole
ffffc90000000000	-22 TB	ffffc97fffffffff	1 TB	virtual memory map (vmemmap_base)
ffffc90000000000	-21 TB	ffffc97fffffffff	1 TB	... unused hole
ffffc90000000000	-20 TB	ffffc97fffffffff	16 TB	KASAN shadow memory
Identical layout to the 56-bit one from here on:				
fffffc0000000000	-4 TB	fffffd7fffffffff	2 TB	... unused hole
fffffc0000000000	-2 TB	fffffe7fffffffff	0.5 TB	vaddr_end for KASLR
fffffc0000000000	-1.5 TB	fffffe7fffffffff	0.5 TB	cpu_entry_area mapping
fffffc0000000000	-1 TB	fffffe7fffffffff	0.5 TB	... unused hole
fffffc0000000000	-512 GB	fffffe7fffffffff	444 GB	%esp fixup stacks
fffffc0000000000	-68 GB	fffffe7fffffffff	64 GB	... unused hole
fffffc0000000000	-4 GB	fffffe7fffffffff	2 GB	EFI region mapping space
fffffc0000000000	-2 GB	fffffe7fffffffff	512 MB	... unused hole
fffffc0000000000	-2048 MB	fffffe7fffffffff	512 MB	kernel text mapping, mapped to physical address 0
fffffc0000000000	-1536 MB	fffffe7fffffffff	1520 MB	module mapping space
fffffc0000000000	-16 MB	fffffe7fffffffff	1520 MB	module mapping space
fffffc0000000000	-11 MB	fffffe7fffffffff	~0.5 MB	kernel-internal fixmap range, variable size and offset
fffffc0000000000	-10 MB	fffffe7fffffffff	4 kB	legacy vsyscall ABI
fffffc0000000000	-2 MB	fffffe7fffffffff	2 MB	... unused hole

Note that negative addresses such as -23TB are absolute addresses in bytes counted down from the top of the 64-bit address space. This makes it slightly easier to visualise the huge 64-bit address space. We can see that user-space virtual memory lies between 0000000000000000 and 00007fffffffffff and kernel space memory lies above ffff800000000000, whence it is clear how the last 16 bits are zero for user-space and 1 for kernel-space: there is a big hole between 0000800000000000 and ffff7fffffffffff which is completely unused. Addresses whose bits 48:63 are either all zero or all one are known as **canonical** addresses, and addresses lying in the hole are called **non-canonical** addresses.

The page tables for kernel-space addresses are setup such that there is a 1:1 mapping between kernel virtual addresses and physical addresses, for easy access of physical memory. The macro `PAGE_OFFSET` yields the address where kernel space begins, and it is this value that must be added to any physical address to convert it to a kernel virtual address. Another important macro is `PAGE_SIZE`, which yields the size of a page and is usually 4096 bytes or 4 KiB for x64 platforms. `PAGE_SHIFT` gives the number of bits 1 must be shifted left to yield `PAGE_SIZE`, and can thus be defined as the logarithm to the base 2 of `PAGE_SIZE`.

2 Page Tables

Let us now look at how exactly virtual addresses are converted into physical ones through page tables, considering specifically the case of four-level page tables. The four different tables are known as

- Page Global Directory
- Page Upper Directory
- Page Middle Directory
- Page Table

Every virtual address (be it kernel or user) is divided into five parts, as shown in figure 1:

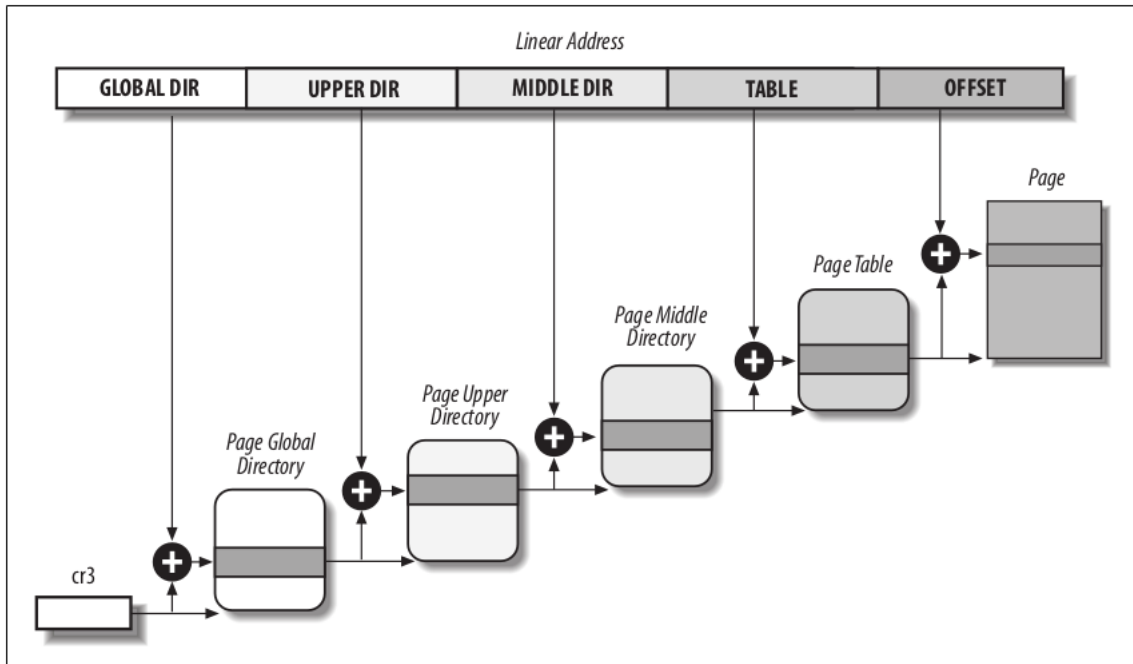


Figure 1: The Linux paging model for four-level page tables

Here, the **OFFSET** part of the address consists of bits 0:11, the **TABLE** part bits 12:20 and so on. The **GLOBAL DIR** part consists of bits 39:47. When a process is created, the kernel creates its corresponding page tables and sets the register `cr3` to point to the base of the Page Global Directory. The **GLOBAL DIR** field within the linear address determines the entry in the Page Global Directory that points to the proper Page Upper Directory. Similarly, the **UPPER DIR** field within the linear address determines the entry in the Page Middle Directory that points to the proper Page Middle Directory. Finally, the Page Table entry gives the physical address of the page itself, and the last 12 bits of the linear address determine the exact address within the page. This process is illustrated in figure 1. A useful macro for calculating the offset is `PAGE_MASK`, which is used to mask all the bits of the **OFFSET** field. This means that `~PAGE_MASK` can be used to calculate the offset from a given virtual address.

Let us now quickly go over how things change for a five-level page table configuration, the one we have to deal with if we're working with kernels later than 4.11. Now, bits 0:56 of the 64 bit address space are used, and the extra page table is simply called `p4d`. Correspondingly, there would be one more table and one more level of indirection in figure 1.

Even more configurations for different 64 bit architectures are given in Table 1

Platform Name	Page Size	Number of address bits used	Number of paging levels	Linear address splitting
alpha	8KiB	43	3	10 + 10 + 10 + 13
ia64	4KiB	39	3	9 + 9 + 9 + 12
ppc64	4KiB	41	3	10 + 10 + 9 + 12
sh64	4KiB	41	3	10 + 10 + 9 + 12
x86-64 [4 level]	4KiB	48	4	9 + 9 + 9 + 9 + 12
x86-64 [5 level]	4KiB	57	5	9 + 9 + 9 + 9 + 9 + 12

Table 1: Paging levels in some 64 bit architectures

3 Our Character Driver - VTP

We are now in a position to start writing our character driver. Our driver, called `vtp` will do the following:

- Accept, through `ioctl`, a user-space virtual address.
- Walk through the `current` process' page tables to find the corresponding physical address
- Print the physical address of the page frame, and the physical address itself to the kernel log buffer

We shall allocate a device number for our char device dynamically using `alloc_chrdev_region`. We can then write a small shell script to find the allocated number through `/proc/devices` and create the corresponding file using `mknod`. The `init` function for our driver is shown in listing 1.

```

1 #include <linux/kernel.h>      /* printk() */
2 #include <linux/slab.h>        /* kmalloc() */
3 #include <linux/fs.h>          /* everything... */
4 #include <linux/errno.h>       /* error codes */
5 #include <linux/cdev.h>        /* char device registration*/
6 #include <linux/types.h>
7 #include <linux/fcntl.h>       /* O_ACCMODE */
8 #include <linux/mm_types.h>     /* struct page and struct mm_struct*/
9 #include <asm/page.h>          /*pgd_t, pte_t, __va etc*/
10 #include <linux/pgtable.h>     /*pgd_offset, pud_offset etc*/
11 #include <asm/pgtable_types.h> /*PTE_PFN_MASK*/
12 #include <linux/highmem.h>
13 #include <asm/io.h>
14
15 MODULE_LICENSE("Dual BSD/GPL");
16
17 int vtp_major = 0;              /*we shall allocate dynamically*/
18 int vtp_minor = 0;
19 int nr_vtp_devs = 1;           /*just one device*/
20
21 struct cdev *vtp_cdev;         /*our character device*/
22 const struct file_operations vtp_fops = {
23     .owner = THIS_MODULE,
24     .unlocked_ioctl = vtp_ioctl, /*we'll only use ioctl*/
25 };
26 int vtp_init(void)
27 {
28     int result, err;
29     dev_t dev = 0;
30
31     result = alloc_chrdev_region(&dev, vtp_minor, nr_vtp_devs,
32                                 "vtp");
33     vtp_major = MAJOR(dev);
34
35     if (result < 0){
36         printk(KERN_WARNING "vtp: can't allocate device number\n");
37         return result;
38     }
39
40     vtp_cdev = cdev_alloc();
41     vtp_cdev->ops = &vtp_fops;
42     vtp_cdev->owner = THIS_MODULE;
43     err = cdev_add(vtp_cdev, dev, 1);
44     if (err)
45         printk(KERN_NOTICE "couldn't add cdev: error %d", err);
46     return 0;
47 }

```

Listing 1: The init function

When writing the script to create the char device, we have to be careful since the script has been executed by **superuser**. Here, we choose to give the group **sudo** read and write access to the device. The script to load the module and create the char device is shown in listing 2.

```

1  #!/bin/sh
2
3  module="vtp"
4  device="vtp"
5  mode="664"          #user, group [rw] others [r]
6
7  #insert module
8  /sbin/insmod ./module.ko $* || exit 1
9
10 #remove devices if already exist
11 rm -f /dev/${device}[0]
12
13 #find the major number
14 major=$(awk "\$2=="$module" {print \$1}" /proc/devices)
15
16 mknod /dev/${device}0 c $major 0
17
18 #now give appropriate permissions to other users, since we have been
19 #invoked by superuser
20 group="sudo"
21
22 chgrp $group /dev/${device}0
23 chmod $mode /dev/${device}0

```

Listing 2: The vtp_load script

Let us now move on to our **ioctl** implementation, where we'll actually walk the page tables. The various page table entries are represented in the kernel by the datatypes **pgd_t**, **p4d_t**, **pud_t**, **pmd_t** and **pte_t**. The address to the Page Global Directory of the current process is stored in the **struct mm_struct** of the **current** process. We can get the correct PGD entry corresponding to our address using the macro **pgd_offset** defined in **linux/pgtable.h**. This macro takes a **struct mm_struct** and a virtual address as a parameter and returns a pointer to the corresponding PGD entry.

There are four other macros, **p4d_offset**, **pud_offset**, **pmd_offset** and **pte_offset_kernel** which serve the same purpose as that of **pgd_offset**. These macros are applied sequentially one after the other, the result of one being the parameter to the next, essentially walking through the page tables. After we have applied **pte_offset_kernel**, what we end up with is the physical address of the page frame corresponding to our virtual address. What remains is to add the last 12 bits of our virtual address (the offset part), and we will have obtained the exact physical address. One last thing to take care of is the fact that **pte_offset_kernel** will give us an address containing various bits for flags such as access rights, page size etc. So, we will have to use **PTE_PFN_MASK**, defined in **asm/pgtable_types.h** to mask off those bits and obtain the actual physical address. The code for our **ioctl** implementation is shown in listing 3.

```

1  long vtp_ioctl(struct file *filp, unsigned int cmd, unsigned long arg)
2  {
3      char c;
4
5      pgd_t *pgd;
6      pte_t *ptep;
7      pud_t *pud;
8      p4d_t *p4d;
9      pmd_t *pmd;
10     char *addr, *pf_addr;
11
12     struct page *page = NULL;
13     struct mm_struct *mm = current->mm;
14
15     pgd = pgd_offset(mm, arg);
16     if (pgd_none(*pgd) || pgd_bad(*pgd))
17         goto out;
18     printk(KERN_NOTICE "Valid pgd\n");
19
20     p4d = p4d_offset(pgd, arg);
21     if (p4d_none(*p4d) || p4d_bad(*p4d))
22         goto out;
23     printk(KERN_NOTICE "Valid p4d\n");
24
25     pud = pud_offset(p4d, arg);
26     if (pud_none(*pud) || pud_bad(*pud))

```

```

27         goto out;
28     printk(KERN_NOTICE "Valid pud\n");
29
30     pmd = pmd_offset(pud, arg);
31     if (pmd_none(*pmd) || pmd_bad(*pmd))
32         goto out;
33     printk(KERN_NOTICE "Valid pmd\n");
34
35     ptep = pte_offset_kernel(pmd, arg);
36     if(!ptep)
37         goto out;
38
39     page = pte_page(*ptep);
40     pf_addr = (char *)((unsigned long)pte_val(*ptep) & PTE_PFN_MASK);
41     addr = pf_addr + (arg & ~PAGE_MASK);
42     c = *((char *) __va(addr));
43     printk(KERN_INFO "the physical address is 0x%px\n", (void *)addr);
44     printk(KERN_INFO "the physical page frame address is 0x%px\n", (void *)pf_addr);
45     printk(KERN_INFO "and the kernel virt addr is 0x%px\n", (void *)__va(addr));
46     printk(KERN_INFO "the byte there is 0x%x\n", c);
47     pte_unmap(ptep);
48
49     return 0;
50
51 out:
52     printk(KERN_INFO "couldn't walk page tables\n");
53     return -1;
54 }

```

Listing 3: Our ioctl function

Notice that we have used %px instead of the familiar %p as a format specifier for pointers. This is because `printk` hashes arguments specified using %p to avoid leaking kernel addresses to user-space. We may write a small program such as the following to test our driver:

```

1 #include <sys/types.h>
2 #include <sys/stat.h>
3 #include <fcntl.h>
4 #include <sys/ioctl.h>
5 #include <stdio.h>
6
7 int main()
8 {
9     char q = 'A';
10    char p = 'B';
11    char *t = malloc(1);
12    int fd = open("/dev/vtp0", O_RDONLY);
13
14    *t='Q';
15
16    ioctl(fd, 0, &q);
17    ioctl(fd, 0, &p);
18    ioctl(fd, 0, t);
19 }

```

If we load our module, run our test program and check the kernel log buffer using `dmesg`, sure enough, everything works perfectly:

```

$ sudo sh vtp_load.sh
$ ./test
$ dmesg

[ 6278.829793] Valid pgd
[ 6278.829795] Valid p4d
[ 6278.829796] Valid pud
[ 6278.829798] Valid pmd
[ 6278.829801] the physical address is 0x0000000125e2be7a
[ 6278.829802] the physical page frame address is 0x0000000125e2b000
[ 6278.829804] and the kernel virt addr is 0xffff9248e5e2be7a
[ 6278.829806] the byte there is 0x41
[ 6278.829809] Valid pgd
[ 6278.829810] Valid p4d
[ 6278.829811] Valid pud
[ 6278.829812] Valid pmd
[ 6278.829813] the physical address is 0x0000000125e2be7b

```

```
[ 6278.829815] the physical page frame address is 0x0000000125e2b000
[ 6278.829816] and the kernel virt addr is 0xffff9248e5e2be7b
[ 6278.829817] the byte there is 0x42
[ 6278.829820] Valid pgd
[ 6278.829821] Valid p4d
[ 6278.829822] Valid pud
[ 6278.829823] Valid pmd
[ 6278.829825] the physical address is 0x000000010fa692a0
[ 6278.829826] the physical page frame address is 0x000000010fa69000
[ 6278.829828] and the kernel virt addr is 0xffff9248cfa692a0
[ 6278.829829] the byte there is 0x51
```