# MIT 6.828: Operating System Engineering

Lab Report

Lab 2: Memory management

CSNLP

csnlp16@126.com

https://csnlp.github.io/

# Contents

# 1   Introduction

Two components of memory management:

1. Physical memory allocator: allocates memory and then free memory.

2. Virtual memory: maps the virtual addresses to physical addresses.

# 2   Part 1: Physical Page Management

The OS must keep track of which parts of physical RAM are free and which are currently in use. In JOS, it keep tracks of which pages are free using a linked list **struct PageInfo** objects, each of which corresponds to a physical page.

## 2.1   Exercise 1

*In the file **kern/pmap.c**, you must implement code for the following functions (probably in the order given).*

1. ***boot_alloc()***,

2. ***mem_init()*** *(only up to the call to check_page_free_list(1)),*

3. ***page_init()***,

4. ***page_alloc()***,

5. ***page_free()***,

*After your implementation, **check_page_free_list()** and **check_page_alloc()** test your physical page allocator. You should boot JOS and see whether check_page_alloc() reports success. Fix your code so that it passes. You may find it helpful to add your own assert()s to verify that your assumptions are correct.*

### 2.1.1   The Analysis of mem_init() Function

Firstly, let's see the **mem_init()** function. The **mem_init()** is used for detect how many memory can be used.

```
1  // Set up a two-level page table:
2  //    kern_pgdir is its linear (virtual) address of the root
3  //
4  // This function only sets up the kernel part of the address space
5  // (ie. addresses >= UTOP).  The user part of the address space
6  // will be set up later.
7  //
8  // From UTOP to ULIM, the user is allowed to read but not write.
9  // Above ULIM the user cannot read or write.
10 void
11 mem_init(void)
12 {
13   uint32_t cr0;
14   size_t n;
15
```

```
16    // Find out how much memory the machine has (npages & npages_basemem).
17    i386_detect_memory();
18
19    // Remove this line when you're ready to test this function.
20    panic("mem_init: This function is not finished\n");
21
22    //////////////////////////////////////////////////////////////////////
23    // create initial page directory.
24    kern_pgdir = (pde_t *) boot_alloc(PGSIZE);
25    memset(kern_pgdir, 0, PGSIZE);
26
27    //////////////////////////////////////////////////////////////////////
28    // Recursively insert PD in itself as a page table, to form
29    // a virtual page table at virtual address UVPT.
30    // (For now, you don't have understand the greater purpose of the
31    // following line.)
32
33    // Permissions: kernel R, user R
34    kern_pgdir[PDX(UVPT)] = PADDR(kern_pgdir) | PTE_U | PTE_P;
35
36    //////////////////////////////////////////////////////////////////////
37    // Allocate an array of npages 'struct PageInfo's and store it in 'pages'.
38    // The kernel uses this array to keep track of physical pages: for
39    // each physical page, there is a corresponding struct PageInfo in this
40    // array.  'npages' is the number of physical pages in memory.  Use memset
41    // to initialize all fields of each struct PageInfo to 0.
42    // Your code goes here:
43
44
45    //////////////////////////////////////////////////////////////////////
46    // Now that we've allocated the initial kernel data structures, we set
47    // up the list of free physical pages. Once we've done so, all further
48    // memory management will go through the page_* functions. In
49    // particular, we can now map memory using boot_map_region
50    // or page_insert
51    page_init();
52
53    check_page_free_list(1);
54    check_page_alloc();
55    check_page();
56
57    //////////////////////////////////////////////////////////////////////
58    // Now we set up virtual memory
59
60    //////////////////////////////////////////////////////////////////////
61    // Map 'pages' read-only by the user at linear address UPAGES
62    // Permissions:
63    //    - the new image at UPAGES -- kernel R, user R
64    //       (ie. perm = PTE_U | PTE_P)
65    //    - pages itself -- kernel RW, user NONE
66    // Your code goes here:
67
68    //////////////////////////////////////////////////////////////////////
69    // Use the physical memory that 'bootstack' refers to as the kernel
70    // stack.  The kernel stack grows down from virtual address KSTACKTOP.
71    // We consider the entire range from [KSTACKTOP-PTSIZE, KSTACKTOP)
72    // to be the kernel stack, but break this into two pieces:
73    //      * [KSTACKTOP-KSTKSIZE, KSTACKTOP) -- backed by physical memory
```

```
74    //     * [KSTACKTOP-PTSIZE, KSTACKTOP-KSTKSIZE) -- not backed; so if
75    //        the kernel overflows its stack, it will fault rather than
76    //        overwrite memory.  Known as a "guard page".
77    //     Permissions: kernel RW, user NONE
78    // Your code goes here:
79
80    //////////////////////////////////////////////////////////////////////
81    // Map all of physical memory at KERNBASE.
82    // Ie.  the VA range [KERNBASE, 2^32) should map to
83    //        the PA range [0, 2^32 - KERNBASE)
84    // We might not have 2^32 - KERNBASE bytes of physical memory, but
85    // we just set up the mapping anyway.
86    // Permissions: kernel RW, user NONE
87    // Your code goes here:
88
89    // Check that the initial page directory has been set up correctly.
90    check_kern_pgdir();
91
92    // Switch from the minimal entry page directory to the full kern_pgdir
93    // page table we just created.  Our instruction pointer should be
94    // somewhere between KERNBASE and KERNBASE+4MB right now, which is
95    // mapped the same way by both page tables.
96    //
97    // If the machine reboots at this point, you've probably set up your
98    // kern_pgdir wrong.
99    lcr3(PADDR(kern_pgdir));
100
101    check_page_free_list(0);
102
103    // entry.S set the really important flags in cr0 (including enabling
104    // paging).  Here we configure the rest of the flags that we care about.
105    cr0 = rcr0();
106    cr0 |= CR0_PE|CR0_PG|CR0_AM|CR0_WP|CR0_NE|CR0_MP;
107    cr0 &= ~(CR0_TS|CR0_EM);
108    lcr0(cr0);
109
110    // Some more checks, only possible after kern_pgdir is installed.
111    check_page_installed_pgdir();
112 }
```

Let's detail this process:

1. We can see that it calls the **i386_detect_memory()** function. Let's check this
   detect function's code:

```
1  static void
2  i386_detect_memory(void)
3  {
4      size_t basemem, extmem, ext16mem, totalmem;
5
6      // Use CMOS calls to measure available base & extended memory.
7      // (CMOS calls return results in kilobytes.)
8      basemem = nvram_read(NVRAM_BASELO);
9      extmem = nvram_read(NVRAM_EXTLO);
10     ext16mem = nvram_read(NVRAM_EXT16LO) * 64;
11
12     // Calculate the number of physical pages available in both base
13     // and extended memory.
14     if (ext16mem)
```

```
15      totalmem = 16 * 1024 + ext16mem;
16    else if (extmem)
17      totalmem = 1 * 1024 + extmem;
18    else
19      totalmem = basemem;
20
21    npages = totalmem / (PGSIZE / 1024);
22    npages_basemem = basemem / (PGSIZE / 1024);
23
24    cprintf("Physical memory: %uK available, base = %uK, extended = %uK\n",
25      totalmem, basemem, totalmem - basemem);
26 }
```

2. After counting the **total memory** and **base memory**, we come to the command as

```
1  /////////////////////////////////////////////////////////////////////
2  // create initial page directory.
3  kern_pgdir = (pde_t *) boot_alloc(PGSIZE);
4  memset(kern_pgdir, 0, PGSIZE);
```

3. We can see that we use **boot_alloc()** function. Let's check this function:

```
1  // This simple physical memory allocator is used only while JOS is setting
2  // up its virtual memory system.  page_alloc() is the real allocator.
3  //
4  // If n>0, allocates enough pages of contiguous physical memory to hold 'n'
5  // bytes.  Doesn't initialize the memory.  Returns a kernel virtual address.
6  //
7  // If n==0, returns the address of the next free page without allocating
8  // anything.
9  //
10 // If we're out of memory, boot_alloc should panic.
11 // This function may ONLY be used during initialization,
12 // before the page_free_list list has been set up.
13 static void *
14 boot_alloc(uint32_t n)
15 {
16        static char *nextfree;  // virtual address of next byte of free memory
17        char *result;
18
19        // Initialize nextfree if this is the first time.
20        // 'end' is a magic symbol automatically generated by the linker,
21        // which points to the end of the kernel's bss segment:
22        // the first virtual address that the linker did *not* assign
23        // to any kernel code or global variables.
24        if (!nextfree) {
25                extern char end[];
26                nextfree = ROUNDUP((char *) end, PGSIZE);
27        }
28
29        // Allocate a chunk large enough to hold 'n' bytes, then update
30        // nextfree.  Make sure nextfree is kept aligned
31        // to a multiple of PGSIZE.
```

```
32          //
33          // LAB 2: Your code here.
34          result = nextfree;
35          nextfree = ROUNDUP(next_free + n, PGSIZE);
36
37          // decide if there is enough nextfree. npages * PGSIZE is
38          // is all available space.
39          if ((uint32_t)nextfree - KERNBASE > (npages * PGSIZE))
40              panic("NO ENOUGH MEMORY!!!\n");
41          return result;
42  }
```

Listing 3: boot_alloc function

It involves the **ROUNDUP**.

```
1  // Round up to the nearest multiple of n
2  #define ROUNDUP(a, n)                                           \
3  ({                                                              \
4          uint32_t __n = (uint32_t) (n);                          \
5          (typeof(a)) (ROUNDDOWN((uint32_t) (a) + __n - 1, __n)); \
6  })
```

We can see that the **boot_alloc** allocate a page of size $n$ which just follows the kernal.

Please notice that **memset()** is defined at lib/string.c.

```
1  #if ASM
2  void *
3  memset(void *v, int c, size_t n)
4  {
5          char *p;
6
7          if (n == 0)
8                  return v;
9          if ((int)v%4 == 0 && n%4 == 0) {
10                 c &= 0xFF;
11                 c = (c<<24)|(c<<16)|(c<<8)|c;
12                 asm volatile("cld; rep stosl\n"
13                         :: "D" (v), "a" (c), "c" (n/4)
14                         : "cc", "memory");
15         } else
16                 asm volatile("cld; rep stosb\n"
17                         :: "D" (v), "a" (c), "c" (n)
18                         : "cc", "memory");
19         return v;
20  }
```

```
1          ///////////////////////////////////////////////////////////////////
2          // Recursively insert PD in itself as a page table, to form
3          // a virtual page table at virtual address UVPT.
4          // (For now, you don't have understand the greater purpose of the
5          // following line.)
6
7          // Permissions: kernel R, user R
8          kern_pgdir[PDX(UVPT)] = PADDR(kern_pgdir) | PTE_U | PTE_P;
```

Listing 4: UVPT

5. Set up the list of physical memory pages:

```
1    //////////////////////////////////////////////////////////////////////
2    // Allocate an array of npages 'struct PageInfo's and store it in 'pages
   ↪ '.
3    // The kernel uses this array to keep track of physical pages: for
4    // each physical page, there is a corresponding struct PageInfo in this
5    // array.  'npages' is the number of physical pages in memory.  Use
   ↪ memset
6    // to initialize all fields of each struct PageInfo to 0.
7    // Your code goes here:
8            // ***
9            //My Code follows
10           // ***
11   pages = (struct PageInfo *) boot_alloc(npages * sizeof(struct PageInfo))
   ↪ ;
12   memset(pages, 0, npages * sizeof(struct PageInfo));
```

Listing 5: My code about the set up of list of physical memory pages

6. The **page_init()** function: see in Section **??** .

## 2.1.2 The Analysis of page_init() Function

```
1  //
2  // Initialize page structure and memory free list.
3  // After this is done, NEVER use boot_alloc again.  ONLY use the page
4  // allocator functions below to allocate and deallocate physical
5  // memory via the page_free_list.
6  //
7  void
8  page_init(void)
9  {
10         // The example code here marks all physical pages as free.
11         // However this is not truly the case.  What memory is free?
12         //  1) Mark physical page 0 as in use.
13         //     This way we preserve the real-mode IDT and BIOS structures
14         //     in case we ever need them.  (Currently we don't, but...)
15         //  2) The rest of base memory, [PGSIZE, npages_basemem * PGSIZE)
16         //     is free.
17         //  3) Then comes the IO hole [IOPHYSMEM, EXTPHYSMEM), which must
18         //     never be allocated.
19         //  4) Then extended memory [EXTPHYSMEM, ...).
20         //     Some of it is in use, some is free. Where is the kernel
21         //     in physical memory?  Which pages are already in use for
22         //     page tables and other data structures?
23         //
24         // Change the code to reflect this.
25         // NB: DO NOT actually touch the physical memory corresponding to
26         // free pages!
27         size_t i;
28         for (i = 0; i < npages; i++) {
29                 pages[i].pp_ref = 0;
30                 pages[i].pp_link = page_free_list;
```

7

```
31          page_free_list = &pages[i];
32      }
33 }
```

We can have our finished **page_init()** function:

```
1  //
2  // Initialize page structure and memory free list.
3  // After this is done, NEVER use boot_alloc again.  ONLY use the page
4  // allocator functions below to allocate and deallocate physical
5  // memory via the page_free_list.
6  //
7  void
8  page_init(void)
9  {
10   // The example code here marks all physical pages as free.
11   // However this is not truly the case.  What memory is free?
12   //  1) Mark physical page 0 as in use.
13   //     This way we preserve the real-mode IDT and BIOS structures
14   //     in case we ever need them.  (Currently we don't, but...)
15   //  2) The rest of base memory, [PGSIZE, npages_basemem * PGSIZE)
16   //     is free.
17   //  3) Then comes the IO hole [IOPHYSMEM, EXTPHYSMEM), which must
18   //     never be allocated.
19   //  4) Then extended memory [EXTPHYSMEM, ...).
20   //     Some of it is in use, some is free. Where is the kernel
21   //     in physical memory?  Which pages are already in use for
22   //     page tables and other data structures?
23   //
24   // Change the code to reflect this.
25   // NB: DO NOT actually touch the physical memory corresponding to
26   // free pages!
27   // Let's begin our implementation:
28
29   // (1). Mark physical page 0 as in use
30   pages[0].pp_ref = 1;
31
32   // (2). The rest of base memory, i.e., the memory between [PGSIZE, npages_basemem*
       ↪ PGSIZE)
33   size_t i;
34   for(i = 1; i < npages_basemem; i++) {
35     pages[i].pp_ref = 0;
36     pages[i].pp_link = page_free_list;
37     page_free_list = &pages[i];
38   }
39
40   // (3). IO hole [IOPHYSMEM, EXTPHYSMEM) must never be allocated.
41
42   for(i=npages_basemem; i < EXTPHYSMEM/PGSIZE; i++) {
43     pages[i].pp_ref = 1;
44   }
45
46   // (4). Extended memory: [EXTPHYSMEM, ...]
47   // The followings are memory that being used in extended memory.
48   physaddr_t first_free_addr = PADDR(boot_alloc(0));
49   size_t first_free_page = first_free_addr/PGSIZE;
50   for(i=EXTPHYSMEM/PGSIZE; i < first_free_page; i++) {
```

```
51    pages[i].pp_ref = 1;
52  }
53  // The followings are memory that are free in extended memory.
54  for(i=first_free_page; i < npages; i++) {
55    pages[i].pp_ref = 0;
56    pages[i].pp_link = page_free_list;
57    page_free_list = &pages[i];
58  }
```

### 2.1.3   Continue of mem_init() Function

After the **page_init()** function, the **mem_init()** function continues:

```
1        check_page_free_list(1);
2        check_page_alloc();
3        check_page();
```

1. **check_page_free_list()** is to check that the pages on **page_free_list** are reasonable.

2. **check_page_alloc()** is to check the physical page allocator, **page_alloc(), page_free(), and page_init()**. Since we have implemented **page_init()**, let's implement the rest two functions.

### 2.1.4   The Analysis of page_alloc() Function

Let's firstly check its description about **page_alloc()**

```
1  //
2  // Allocates a physical page.  If (alloc_flags & ALLOC_ZERO), fills the entire
3  // returned physical page with '\0' bytes.  Does NOT increment the reference
4  // count of the page - the caller must do these if necessary (either explicitly
5  // or via page_insert).
6  //
7  // Be sure to set the pp_link field of the allocated page to NULL so
8  // page_free can check for double-free bugs.
9  //
10 // Returns NULL if out of free memory.
11 //
12 // Hint: use page2kva and memset
13 struct PageInfo *
14 page_alloc(int alloc_flags)
15 {
16        // Fill this function in
17        return 0;
18 }
```

We notice the hint the **page2kva, memset**, defined in **kern/pmap.h** may be useful: let see this two function:

```
1  static inline void*
2  page2kva(struct PageInfo *pp)
3  {
4        return KADDR(page2pa(pp));
5  }
```

Listing 7: page2kva inline

.

- **KADDR** is about:

```
/* This macro takes a physical address and returns the corresponding kernel
 * virtual address.  It panics if you pass an invalid physical address. */
#define KADDR(pa) _kaddr(__FILE__, __LINE__, pa)

static inline void*
_kaddr(const char *file, int line, physaddr_t pa)
{
        if (PGNUM(pa) >= npages)
                _panic(file, line, "KADDR called with invalid pa %08lx", pa);
        return (void *)(pa + KERNBASE);
}
```

- **page2pa** is for return the physical address given **PageInfo \***

```
static inline physaddr_t
page2pa(struct PageInfo *pp)
{
        return (pp - pages) << PGSHIFT;
}
```

Let's see the final implementation:

```
struct PageInfo *
page_alloc(int alloc_flags)
{
        // Fill this function in
        struct PageInfo *result;

        // Out of free memory
        if(page_free_list == NULL) {
                return NULL;
        }
        // if have free memory, return the Page that page_free_list points at.
        // move the page_free_list to next item in the linked list.
        result = page_free_list;
        page_free_list = result->pp_link;
        result->pp_link = NULL;

        // alloc_flags & ALLOC_ZERO
        if (alloc_flags & ALLOC_ZERO) {
                memset(page2kva(result), 0, PGSIZE);
        }
        return result;
        //return 0;
}
```

### 2.1.5   The Analysis of page_free Function

Finally, it's the **page_free** function: Let's see the original code and its hints:

```c
//
// Return a page to the free list.
// (This function should only be called when pp->pp_ref reaches 0.)
//
void
page_free(struct PageInfo *pp)
{
        // Fill this function in
        // Hint: You may want to panic if pp->pp_ref is nonzero or
        // pp->pp_link is not NULL.
}
```

We have two steps to do:

1. double check **pp → pp_link** and **pp → pp_ref**.

2. insert the *freed* page to **free_page_list**.

Now, we have our finished code

```c
void
page_free(struct PageInfo *pp)
{
        // Fill this function in
        // Hint: You may want to panic if pp->pp_ref is nonzero or
        // pp->pp_link is not NULL.
        if(pp->pp_link || pp->pp_ref) {
                panic("pp->pp_ref is nonzero or pp->pp_link is not NULL");
        }
        pp->pp_link = page_free_list;
        page_free_list = pp;
}
```

Then, we finally finish this part.

.

```
Selector   +-------------+          +----------+
---------->|             |          |          |
           | Segmentation |          | Paging   |
Software   |             |--------->|          |---------> RAM
    Offset | Mechanism   |          | Mechanism|
---------->|             |          |          |
           +-------------+          +----------+
    Virtual              Linear              Physical
```
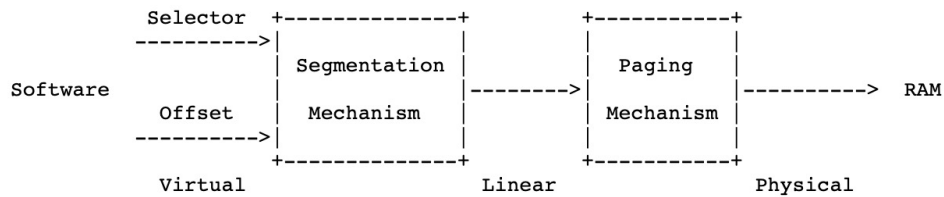
Figure 1: The relationships of virtual address, linear address, and physical address.

# 3   Part 2: Virtual Memory

## 3.1   x86's Protected-mode Memory Management Architecture

### 3.1.1   Segmentation

### 3.1.2   Page Translation

## 3.2   Exercise 2

*Look at chapters 5 and 6 of the Intel 80386 Reference Manual, if you haven't done so already. Read the sections about page translation and page-based protection closely (5.2 and 6.4). We recommend that you also skim the sections about segmentation; while JOS uses the paging hardware for virtual memory and protection, segment translation and segment-based protection cannot be disabled on the x86, so you will need a basic understanding of it.*

## 3.3   The Relationships of Virtual Address, Linear Address, and Physical Address

In x86 terminology, a virtual address consists of a segment selector and an offset within the segment. A linear address is what you get after segment translation but before page translation. A physical address is what you finally get after both segment and page translation and what ultimately goes out on the hardware bus to your RAM.

## 3.4   Exercise 3

## 3.5   Virtual Address, Physical Address, and Their Interconversion

| C type | Address type |
|--------|--------------|
| uintptr_t | Virtual |
| physaddr_t | Physical |

For code executing on the CPU, once we're in the protected mode, there's no way to directly use a linear address or physical address. **All memory reference are interpreted as virtual address and translated by the MMU, which means all pointer in C are virtual addresses.**

Let's see the interconversion:

```
1  /* This macro takes a kernel virtual address -- an address that points above
2   * KERNBASE, where the machine's maximum 256MB of physical memory is mapped --
3   * and returns the corresponding physical address.  It panics if you pass it a
4   * non-kernel virtual address.
5   */
6  #define PADDR(kva) _paddr(__FILE__, __LINE__, kva)
7
8  static inline physaddr_t
9  _paddr(const char *file, int line, void *kva)
10 {
11         if ((uint32_t)kva < KERNBASE)
12                 _panic(file, line, "PADDR called with invalid kva %08lx", kva);
13         return (physaddr_t)kva - KERNBASE;
14 }
15
16 /* This macro takes a physical address and returns the corresponding kernel
17  * virtual address.  It panics if you pass an invalid physical address. */
18 #define KADDR(pa) _kaddr(__FILE__, __LINE__, pa)
19
20 static inline void*
21 _kaddr(const char *file, int line, physaddr_t pa)
22 {
23         if (PGNUM(pa) >= npages)
24                 _panic(file, line, "KADDR called with invalid pa %08lx", pa);
25         return (void *)(pa + KERNBASE);
26 }
```

Listing 8: Interconversion of physical address and virtual address.

## 3.6 Reference Counting

A same physical page can be mapped at multiple virtual addresses simultaneously. The number of references to each physical page is kept in the **pp_ref** field of the **struct PageInfo** corresponding to the physical page. When **pp_ref** goes to zero for a physical page, the page can be freed since it is no longer being used.

## 3.7 Page Table Management

Exercise 4 is about a set of routines to manage page tables: to insert and remove linear-to-physical mappings, to create page table pages when needed and so on.

### 3.7.1 Exercise 4

*In the file kern/pmap.c, you must implement code for the following functions:*

- *pgdir_walk()*

- *boot_map_region()*

- *page_lookup()*

- *page_remove()*

- *page_insert()*

*check_page(), called from mem_init() tests your page table management routines. Your should make sure it reports success before proceeding.*
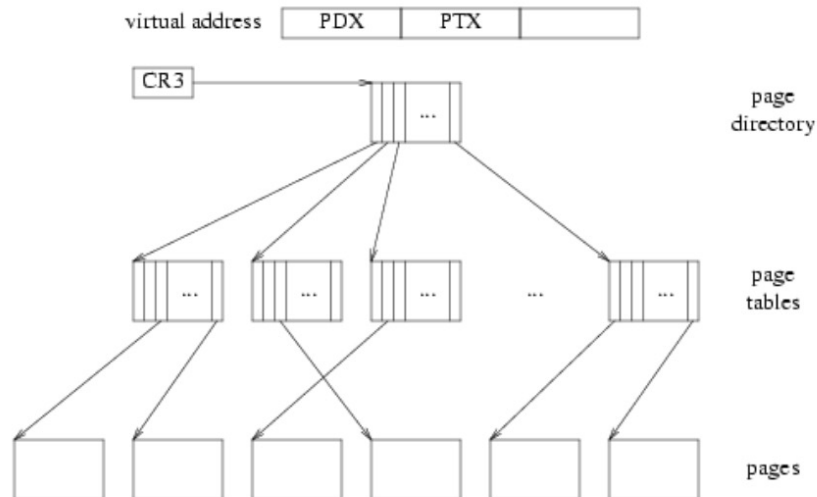
### 3.7.2 Some Useful Function in inc/mmu.h

```
1  /*
2   * This file contains definitions for the x86 memory management unit (MMU),
3   * including paging- and segmentation-related data structures and constants,
4   * the %cr0, %cr4, and %eflags registers, and traps.
5   */
6
7  /*
8   *
9   *  Part 1.  Paging data structures and constants.
10  *
11  */
12
13 // A linear address 'la' has a three-part structure as follows:
14 //
15 // +--------10------+-------10-------+---------12----------+
16 // | Page Directory |   Page Table   | Offset within Page  |
17 // |      Index     |      Index     |                     |
18 // +----------------+----------------+---------------------+
19 //  \--- PDX(la) --/ \--- PTX(la) --/ \---- PGOFF(la) ----/
20 //  \---------- PGNUM(la) ----------/
21 //
22 // The PDX, PTX, PGOFF, and PGNUM macros decompose linear addresses as shown.
23 // To construct a linear address la from PDX(la), PTX(la), and PGOFF(la),
24 // use PGADDR(PDX(la), PTX(la), PGOFF(la)).
25
26 // page number field of address
27 #define PGNUM(la) (((uintptr_t) (la)) >> PTXSHIFT)
28
29 // page directory index
30 #define PDX(la)   ((((uintptr_t) (la)) >> PDXSHIFT) & 0x3FF)
31
32 // page table index
33 #define PTX(la)   ((((uintptr_t) (la)) >> PTXSHIFT) & 0x3FF)
34
35 // offset in page
36 #define PGOFF(la) (((uintptr_t) (la)) & 0xFFF)
37
38 // construct linear address from indexes and offset
39 #define PGADDR(d, t, o) ((void*) ((d) << PDXSHIFT | (t) << PTXSHIFT | (o)))
```

Let's get more information about **PDX** and **PTX**.
1. CR3 points at the page directory.

2. The PDX part of the address indexes into the page directory to give you a page table.

3. The PTX part indexed into the page table to give you a page.

However, the processor has no concept of page directories, page tables and pages. The processor just follows pointers:

```
1  pd = lcr3();
2  pt = *(pd + 4*PDX);
3  page = *(pt + 4*PTX);
```

### 3.7.3  pgdir_walk() Function

Let's firstly see its hints:

```
1  // Given 'pgdir', a pointer to a page directory, pgdir_walk returns
2  // a pointer to the page table entry (PTE) for linear address 'va'.
3  // This requires walking the two-level page table structure.
4  //
5  // The relevant page table page might not exist yet.
6  // If this is true, and create == false, then pgdir_walk returns NULL.
7  // Otherwise, pgdir_walk allocates a new page table page with page_alloc.
8  //    - If the allocation fails, pgdir_walk returns NULL.
9  //    - Otherwise, the new page's reference count is incremented,
10 //      the page is cleared,
11 //      and pgdir_walk returns a pointer into the new page table page.
12 //
13 // Hint 1: you can turn a PageInfo * into the physical address of the
14 // page it refers to with page2pa() from kern/pmap.h.
15 //
16 // Hint 2: the x86 MMU checks permission bits in both the page directory
17 // and the page table, so it's safe to leave permissions in the page
18 // directory more permissive than strictly necessary.
19 //
20 // Hint 3: look at inc/mmu.h for useful macros that manipulate page
21 // table and page directory entries.
22 //
23 pte_t *
24 pgdir_walk(pde_t *pgdir, const void *va, int create)
25 {
26         // Fill this function in
27         return NULL;
28 }
```

See the implementation:

```
pte_t *
pgdir_walk(pde_t *pgdir, const void *va, int create)
{
        // Fill this function in
        assert(pgdir != NULL);
        unsigned int dir_offset = PDX(va);
        unsigned int page_offset = PTX(va);
        pte_t *page_table = NULL;
        pde_t *dir_page_entry = pgdir + dir_offset;

        struct PageInfo *new_page = NULL;

        bool pg_dir_entry_exists_flag = (*dir_page_entry) & PTE_P;
        if (!pg_dir_entry_exists_flag){
                if (!create)
                        return NULL;
                // doesn't exist but want to create
                else {
                        new_page = page_alloc(ALLOC_ZERO);
                        if (new_page == NULL)
                                return NULL;
                        new_page->pp_ref++;
                        *dir_page_entry = (page2pa(new_page)|PTE_P|PTE_W|PTE_U);
                }
        }
        // if exists, just return.
        page_table = KADDR(PTE_ADDR(*dir_page_entry));
        return &page_table[page_offset];
}
```

### 3.7.4  boot_map_region() Function

```
//
// Map [va, va+size) of virtual address space to physical [pa, pa+size)
// in the page table rooted at pgdir.  Size is a multiple of PGSIZE, and
// va and pa are both page-aligned.
// Use permission bits perm|PTE_P for the entries.
//
// This function is only intended to set up the ''static'' mappings
// above UTOP. As such, it should *not* change the pp_ref field on the
// mapped pages.
//
// Hint: the TA solution uses pgdir_walk
static void
boot_map_region(pde_t *pgdir, uintptr_t va, size_t size, physaddr_t pa, int perm)
{
        // Fill this function in

}
```

This function is only intended to set up the static mappings above UTOP.

```
static void
boot_map_region(pde_t *pgdir, uintptr_t va, size_t size, physaddr_t pa, int perm)
{
        // Fill this function in
```

```
5          pte_t *entry = NULL;
6          for (int i=0; i < size; i+=PGSIZE) {
7                  // get the entry corresponding to va.
8                  entry = pgdir_walk(pgdir, (void*)va, 1);
9                  *entry = (pa | perm | PTE_P);
10                 pa += PGSIZE;
11                 va += PGSIZE;
12         }
13 }
```

### 3.7.5  page_insert() Function

```
1  //
2  // Map the physical page 'pp' at virtual address 'va'.
3  // The permissions (the low 12 bits) of the page table entry
4  // should be set to 'perm|PTE_P'.
5  //
6  // Requirements
7  //    - If there is already a page mapped at 'va', it should be page_remove()d.
8  //    - If necessary, on demand, a page table should be allocated and inserted
9  //      into 'pgdir'.
10 //    - pp->pp_ref should be incremented if the insertion succeeds.
11 //    - The TLB must be invalidated if a page was formerly present at 'va'.
12 //
13 // Corner-case hint: Make sure to consider what happens when the same
14 // pp is re-inserted at the same virtual address in the same pgdir.
15 // However, try not to distinguish this case in your code, as this
16 // frequently leads to subtle bugs; there's an elegant way to handle
17 // everything in one code path.
18 //
19 // RETURNS:
20 //    0 on success
21 //    -E_NO_MEM, if page table couldn't be allocated
22 //
23 // Hint: The TA solution is implemented using pgdir_walk, page_remove,
24 // and page2pa.
25 //
26 int
27 page_insert(pde_t *pgdir, struct PageInfo *pp, void *va, int perm)
28 {
29   // Fill this function in
30   return 0;
31 }
```

### 3.7.6  page_lookup() Function

The origin code of **page_lookup()** function

```
1  //
2  // Return the page mapped at virtual address 'va'.
3  // If pte_store is not zero, then we store in it the address
4  // of the pte for this page.  This is used by page_remove and
5  // can be used to verify page permissions for syscall arguments,
6  // but should not be used by most callers.
7  //
8  // Return NULL if there is no page mapped at va.
9  //
10 // Hint: the TA solution uses pgdir_walk and pa2page.
11 //
```

```
12  struct PageInfo *
13  page_lookup(pde_t *pgdir, void *va, pte_t **pte_store)
14  {
15          // Fill this function in
16          return NULL;
17  }
```

```
1   struct PageInfo *
2   page_lookup(pde_t *pgdir, void *va, pte_t **pte_store)
3   {
4           // Fill this function in
5           pte_t *pte_ptr = NULL;
6           struct PageInfo *ret = NULL;
7
8           // find the va's corresponding page entry.
9           pte_ptr = pgdir_walk(pgdir, va, 0);
10
11          // if there is no page mapped at va
12          if (pte_ptr == NULL)
13                  return NULL;
14          if (!(*entry & PTE_P))
15                  return NULL;
16
17          ret = pa2page(PTE_ADDR(*pte_ptr));
18          if (pte_store != 0) {
19                  *pte_store = pte_ptr;
20          }
21          return ret;
22  }
```

### 3.7.7 page_remove Function

```
1   //
2   // Unmaps the physical page at virtual address 'va'.
3   // If there is no physical page at that address, silently does nothing.
4   //
5   // Details:
6   //   - The ref count on the physical page should decrement.
7   //   - The physical page should be freed if the refcount reaches 0.
8   //   - The pg table entry corresponding to 'va' should be set to 0.
9   //     (if such a PTE exists)
10  //   - The TLB must be invalidated if you remove an entry from
11  //     the page table.
12  //
13  // Hint: The TA solution is implemented using page_lookup,
14  //       tlb_invalidate, and page_decref.
15  //
16  void
17  page_remove(pde_t *pgdir, void *va)
18  {
19          // Fill this function in
20  }
```

Now, we have finished part 1 and part 2 in Lab 2.

# 4 Part 3: Kernal Address Space

JOS divides the processor's 32-bit linear address space into two parts:

1. The lower part: usually controlled by the user environment.

2. The upper part: usually controlled by the kernal.

The dividing line is defined somewhat arbitrarily by the symbol *ULIM* in **inc/memlayout.h**

## 4.1 Permissions and Fault Isolation

From low to top, there are mainly three parts:

1. **The address below *UTOP***: for the user environment to use. The user environment will set the permissions for accessing this memory.

2. **The address range [*UTOP, ULIM*)**: both the kernal and the user environment have the same permission: they can read but not write this address range.

3. **The address above *ULIM***: the user environment have no permission to any of this range of memory, while the kernal can be able to read and write this memory.

The user environment will have no permission to any of the memory above *ULIM*, while the kernal have the right to read and write this memory.

## 4.2 Initializing the Kernel Address Space

### 4.2.1 Exercise 5

*Fill in the missing code in mem_init() after the call to check_page().*
*Your code should now pass the check_kern_pgdir() and check_page_installed_pgdir() checks.*
   Firstly see the hints:

```
1       ////////////////////////////////////////////////////////////////////
2       // Now we set up virtual memory
3
4       ////////////////////////////////////////////////////////////////////
5       // Map 'pages' read-only by the user at linear address UPAGES
6       // Permissions:
7       //    - the new image at UPAGES -- kernel R, user R
8       //      (ie. perm = PTE_U | PTE_P)
9       //    - pages itself -- kernel RW, user NONE
10      // Your code goes here:
11
12      ////////////////////////////////////////////////////////////////////
13      // Use the physical memory that 'bootstack' refers to as the kernel
14      // stack.  The kernel stack grows down from virtual address KSTACKTOP.
15      // We consider the entire range from [KSTACKTOP-PTSIZE, KSTACKTOP)
16      // to be the kernel stack, but break this into two pieces:
17      //     * [KSTACKTOP-KSTKSIZE, KSTACKTOP) -- backed by physical memory
18      //     * [KSTACKTOP-PTSIZE, KSTACKTOP-KSTKSIZE) -- not backed; so if
19      //        the kernel overflows its stack, it will fault rather than
```

```
20         //          overwrite memory.   Known as a "guard page".
21         //      Permissions: kernel RW, user NONE
22         // Your code goes here:
23
24         //////////////////////////////////////////////////////////////////////
25         // Map all of physical memory at KERNBASE.
26         // Ie.   the VA range [KERNBASE, 2^32) should map to
27         //       the PA range [0, 2^32 - KERNBASE)
28         // We might not have 2^32 - KERNBASE bytes of physical memory, but
29         // we just set up the mapping anyway.
30         // Permissions: kernel RW, user NONE
31         // Your code goes here:
```

We can see that we have three parts two finish:

1. Map 'pages' read-only by the user at linear address UPAGES

2. Use the physical memory that 'bootstack' refers to as the kernel stack.

3.

# 5 Reference

1. bysui's github and blog

   - github: https://github.com/bysui/mit6.828
   - blog: https://blog.csdn.net/bysui

2. SmallPond's github and blog

   - github: https://github.com/SmallPond/MIT6.828_OS
   - blog: https://me.csdn.net/Small_Pond

3. SimpCosm's github

   - github: https://github.com/SimpCosm/6.828

4. fatsheep9146's blog

   - blog: https://www.cnblogs.com/fatsheep9146/category/769143.html