
**Report for lab2, Kexing Zhou,
1900013008**

Contents

Environment Configuration	2
Test Compiler Toolchain	3
QEMU Emulator	3
Memory Management	3
Exercise 1	3
pmap.c, boot_alloc	3
pmap.c, mem_init	3
pmap.c, page_init	4
pmap.c, page_alloc	5
pmap.c, page_free	5
Exercise 2	5
Exercise 3	5
Question	6
Exercise 4	6
pmap.c, pgdir_walk	6
pmap.c, boot_map_region	6
pmap.c, page_lookup	6
pmap.c, page_remove	7
pmap.c, page_insert	7
Exercise 5	7
Question	8
Question 2	8
Question 3	11
Question 4	11
Question 5	11
Question 6	11
Challenge 1	12
Challenge 2	13

[TOC]

Environment Configuration

```
1 Hardware Environment:
2 Memory:             16GB
3 Processor:          Intel Core i7-8550U CPU @ 1.66GHz 8
4 GPU:                NVIDIA GeForce RTX 2070
5 OS Type:            64 bit
6 Disk:               924GB
7
8 Software Environment:
9 OS:                  Arch Linux
10 Gcc:                 Gcc 11.1.0
11 Make:                GNU Make 4.3
12 Gdb:                 GNU gdb 11.1
```

Test Compiler Toolchain

```
1 $ objdump -i # the 5th line say elf32-i386
2 $ gcc -m32 -print-libgcc-file-name
3 /usr/lib/gcc/x86_64-pc-linux-gnu/11.1.0/32/libgcc.a
```

QEMU Emulator

```
1 $ sudo pacman -S riscv64-linux-gnu-binutils \
2   riscv64-linux-gnu-gcc riscv64-linux-gnu-gdb qemu-arch-extra
```

Memory Management

Exercise 1

pmap.c, boot_alloc

```
1 static void *
2 boot_alloc(uint32_t n)
3 {
4     static char *nextfree; // virtual address of next byte of free memory
5     char *result;
6
7     // Initialize nextfree if this is the first time.
8     // 'end' is a magic symbol automatically generated by the linker,
9     // which points to the end of the kernel's bss segment:
10    // the first virtual address that the linker did *not* assign
11    // to any kernel code or global variables.
12    if (!nextfree) {
13        extern char end[];
14        nextfree = ROUNDUP((char *) end, PGSIZE);
15    }
16
17    // Allocate a chunk large enough to hold 'n' bytes, then update
18    // nextfree. Make sure nextfree is kept aligned
19    // to a multiple of PGSIZE.
20
21    void *ret = nextfree;
22    nextfree = ROUNDUP(nextfree + n, PGSIZE);
23
24    return ret;
25 }
```

pmap.c, mem_init

```
1 void
2 mem_init(void)
3 {
4     uint32_t cr0;
5     size_t n;
6
7     // Find out how much memory the machine has (npages & npages_baseem).
8     i386_detect_memory();
9
10    // Remove this line when you're ready to test this function.
11
12    //////////////////////////////////////
13    // create initial page directory.
14    kern_pgdir = (pde_t *) boot_alloc(PGSIZE);
15    memset(kern_pgdir, 0, PGSIZE);
16
17    //////////////////////////////////////
18    // Recursively insert PD in itself as a page table, to form
19    // a virtual page table at virtual address UVPT.
20    // (For now, you don't have understand the greater purpose of the
21    // following line.)
22
23    // Permissions: kernel R, user R
24    kern_pgdir[PDX(UVPT)] = PADDR(kern_pgdir) | PTE_U | PTE_P;
25
26    //////////////////////////////////////
27    // Allocate an array of npages 'end's and store it in 'pages'.
28    // The kernel uses this array to keep track of physical pages: for
29    // each physical page, there is a corresponding struct PageInfo in this
```

```

30 // array. 'npages' is the number of physical pages in memory. Use memset
31 // to initialize all fields of each struct PageInfo to 0.
32 pages = boot_alloc(sizeof(*pages) * npages);
33 memset(pages, 0, sizeof(*pages) * npages);
34
35 ////////////////////////////////////////////////////
36 // Now that we've allocated the initial kernel data structures, we set
37 // up the list of free physical pages. Once we've done so, all further
38 // memory management will go through the page_* functions. In
39 // particular, we can now map memory using boot_map_region
40 // or page_insert
41 page_init();
42
43 check_page_free_list(1);
44 check_page_alloc();
45 check_page();
46
47 ////////////////////////////////////////////////////
48 // Now we set up virtual memory
49
50 ////////////////////////////////////////////////////
51 // Map 'pages' read-only by the user at linear address UPAGES
52 // Permissions:
53 //   - the new image at UPAGES -- kernel R, user R
54 //   (ie. perm = PTE_U | PTE_P)
55 //   - pages itself -- kernel RW, user NONE
56 boot_map_region(kern_pgdir, UPAGES, PTSIZE, PADDR(pages), PTE_P | PTE_W);
57
58 ////////////////////////////////////////////////////
59 // Use the physical memory that 'bootstack' refers to as the kernel
60 // stack. The kernel stack grows down from virtual address KSTACKTOP.
61 // We consider the entire range from [KSTACKTOP-PTSIZE, KSTACKTOP)
62 // to be the kernel stack, but break this into two pieces:
63 //   * [KSTACKTOP-KSTKSIZE, KSTACKTOP) -- backed by physical memory
64 //   * [KSTACKTOP-PTSIZE, KSTACKTOP-KSTKSIZE) -- not backed; so if
65 //     the kernel overflows its stack, it will fault rather than
66 //     overwrite memory. Known as a "guard page".
67 // Permissions: kernel RW, user NONE
68 boot_map_region(kern_pgdir, KSTACKTOP - KSTKSIZE, KSTKSIZE, PADDR(bootstack), PTE_P | PTE_W);
69
70 ////////////////////////////////////////////////////
71 // Map all of physical memory at KERNBASE.
72 // Ie. the VA range [KERNBASE, 2^32) should map to
73 // the PA range [0, 2^32 - KERNBASE)
74 // We might not have 2^32 - KERNBASE bytes of physical memory, but
75 // we just set up the mapping anyway.
76 // Permissions: kernel RW, user NONE
77 boot_map_region(kern_pgdir, KERNBASE, -KERNBASE, 0, PTE_P | PTE_W);
78
79
80 // Check that the initial page directory has been set up correctly.
81 check_kern_pgdir();
82
83 uint32_t cr4 = rcr4();
84 cr4 |= CR4_PSE;
85 lcr4(cr4);
86
87 // Switch from the minimal entry page directory to the full kern_pgdir
88 // page table we just created. Our instruction pointer should be
89 // somewhere between KERNBASE and KERNBASE+4MB right now, which is
90 // mapped the same way by both page tables.
91 //
92 // If the machine reboots at this point, you've probably set up your
93 // kern_pgdir wrong.
94 lcr3(PADDR(kern_pgdir));
95
96 check_page_free_list(0);
97
98 // entry.S set the really important flags in cr0 (including enabling
99 // paging). Here we configure the rest of the flags that we care about.
100 cr0 = rcr0();
101 cr0 |= CR0_PE|CR0_PG|CR0_AM|CR0_WP|CR0_NE|CR0_MP;
102 cr0 &= ~(CR0_TS|CR0_EM);
103 lcr0(cr0);
104
105 // Some more checks, only possible after kern_pgdir is installed.
106 check_page_installed_pgdir();
107 }

```

pmap.c, page_init

```

1 void
2 page_init(void)
3 {
4     // The example code here marks all physical pages as free.
5     // However this is not truly the case. What memory is free?
6     // 1) Mark physical page 0 as in use.
7     //     This way we preserve the real-mode IDT and BIOS structures
8     //     in case we ever need them. (Currently we don't, but...)
9     // 2) The rest of base memory, [PGSIZE, npages_basemem * PGSIZE)
10    //    is free.
11    // 3) Then comes the IO hole [IOIOPHYSMEMPHYSMEM, EXTPHYSMEM), which must
12    //    never be allocated.

```

```
13 // 4) Then extended memory [EXTPHYSMEM, ...).
14 // Some of it is in use, some is free. Where is the kernel
15 // in physical memory? Which pages are already in use for
16 // page tables and other data structures?
17 //
18 // Change the code to reflect this.
19 // NB: DO NOT actually touch the physical memory corresponding to
20 // free pages!
21
22 page_free_list = NULL;
23
24 for(size_t i = PGNUM(PADDR(boot_alloc(0))); i < npages; i++) {
25     pages[i].pp_ref = 0;
26     pages[i].pp_link = page_free_list;
27     page_free_list = &pages[i];
28 }
29 // at system start, the lower memory is mapped into the initial pagetable
30 // so I put these pages into the top of page_free_list
31 for(size_t i = 1; i < npages_basemem; i++) {
32     pages[i].pp_ref = 0;
33     pages[i].pp_link = page_free_list;
34     page_free_list = &pages[i];
35 }
36 }
```

pmap.c, page_alloc

```
1 struct PageInfo *
2 page_alloc(int alloc_flags)
3 {
4     if(!page_free_list)
5         return NULL;
6     struct PageInfo * ret = page_free_list;
7     page_free_list = ret->pp_link;
8     ret->pp_link = NULL;
9     if(alloc_flags & ALLOC_ZERO) {
10         memset(page2kva(ret), 0, PGSIZE);
11     }
12     return ret;
13 }
```

pmap.c, page_free

```
1 void
2 page_free(struct PageInfo *pp)
3 {
4     // Fill this function in
5     // Hint: You may want to panic if pp->pp_ref is nonzero or
6     // pp->pp_link is not NULL.
7     if(pp->pp_ref || pp->pp_link)
8         panic("page to free is already in free list\n");
9     pp->pp_link = page_free_list;
10     page_free_list = pp;
11 }
```

Exercise 2

Nothing to report.

Exercise 3

Use the `xp` command in the QEMU monitor and the `x` command in GDB to inspect memory at corresponding physical and virtual addresses and make sure you see the same data.

In QEMU

```
1 (qemu) xp 0x100000
2 0000000000100000: 0x1badb002
```

In GDB:

```
1 (gdb) p/x *0xf0100000
2 $1 = 0x1badb002
```

Question

Assuming that the following JOS kernel code is correct, what type should variable `x` have, `uintptr_t` or `physaddr_t`?

It should be `uintptr_t`.

Exercise 4

`pmap.c`, `pgdir_walk`

```
1 pte_t *
2 pgdir_walk(pde_t *pgdir, const void *va, int create) {
3     pde_t pde = pgdir[PDX(va)];
4     if(pde & PTE_P) {
5         pte_t *ptab = KADDR(PTE_ADDR(pde));
6         return &ptab[PTX(va)];
7     }
8     else if(create) {
9         struct PageInfo *ptab_info = page_alloc(ALLOC_ZERO);
10        if(ptab_info == NULL) return NULL;
11        ptab_info->pp_ref++;
12        physaddr_t pa = page2pa(ptab_info);
13        pgdir[PDX(va)] = pa | PTE_P | PTE_U | PTE_W;
14        pte_t *ptab = KADDR(pa);
15        return &ptab[PTX(va)];
16    }
17    else return NULL;
18 }
```

`pmap.c`, `boot_map_region`

```
1 static void
2 boot_map_region(pde_t *pgdir, uintptr_t va, size_t size, physaddr_t pa, int perm)
3 {
4     perm = (perm & 0x3FF) | PTE_P;
5     for(size_t offset = 0; offset < size; offset += PGSIZE) {
6         pte_t *ppte = pgdir_walk(pgdir, (void*)va + offset, true);
7         if(ppte == NULL) panic("No Available Page");
8         *ppte = (pa + offset) | perm;
9     }
10 }
```

`pmap.c`, `page_lookup`

```
1 struct PageInfo *
2 page_lookup(pde_t *pgdir, void *va, pte_t **pte_store)
3 {
4     pte_t *pptab = pgdir_walk(pgdir, va, false);
5     if(pte_store) {
6         *pte_store = pptab;
7     }
8     if(pptab && (*pptab & PTE_P)) {
9         return pa2page(PTE_ADDR(*pptab));
10    }
11    else {
12        return NULL;
13    }
14 }
```

pmap.c, page_remove

```
1 void
2 page_remove(pde_t *pgdir, void *va)
3 {
4     pte_t *ppte;
5     struct PageInfo *info = page_lookup(pgdir, va, &ppte);
6     if(info == NULL) return;
7     *ppte = 0;
8     page_decref(info);
9     tlb_invalidate(pgdir, va);
10 }
```

pmap.c, page_insert

```
1 int
2 page_insert(pde_t *pgdir, struct PageInfo *pp, void *va, int perm)
3 {
4     perm = (perm & 0x3FF) | PTE_P;
5     pte_t *ppte = pgdir_walk(pgdir, va, true);
6     if(ppte == NULL) {
7         return -E_NO_MEM;
8     }
9     physaddr_t pa = page2pa(pp);
10    bool same_map = false;
11    if(*ppte & PTE_P) {
12        if(PTE_ADDR(*ppte) != pa) {
13            page_remove(pgdir, va);
14        }
15        else {
16            same_map = true;
17        }
18    }
19    if(!same_map) pp->pp_ref++;
20    *ppte = pa | perm;
21    return 0;
22 }
```

Exercise 5

Fill in the missing code in mem_init() after the call to check_page().

```
1 void
2 mem_init(void)
3 {
4     uint32_t cr0;
5     size_t n;
6
7     // Find out how much memory the machine has (npages & npages_baseemem).
8     i386_detect_memory();
9
10    // Remove this line when you're ready to test this function.
11
12    //////////////////////////////////////
13    // create initial page directory.
14    kern_pgdir = (pde_t *) boot_alloc(PGSIZE);
15    memset(kern_pgdir, 0, PGSIZE);
16
17    //////////////////////////////////////
18    // Recursively insert PD in itself as a page table, to form
19    // a virtual page table at virtual address UVPT.
20    // (For now, you don't have understand the greater purpose of the
21    // following line.)
22
23    // Permissions: kernel R, user R
24    kern_pgdir[PDX(UVPT)] = PADDR(kern_pgdir) | PTE_U | PTE_P;
25
26    //////////////////////////////////////
27    // Allocate an array of npages 'end's and store it in 'pages'.
28    // The kernel uses this array to keep track of physical pages: for
29    // each physical page, there is a corresponding struct PageInfo in this
30    // array. 'npages' is the number of physical pages in memory. Use memset
31    // to initialize all fields of each struct PageInfo to 0.
32    pages = boot_alloc(sizeof(*pages) * npages);
33    memset(pages, 0, sizeof(*pages) * npages);
34
35    //////////////////////////////////////
36    // Now that we've allocated the initial kernel data structures, we set
37    // up the list of free physical pages. Once we've done so, all further
38    // memory management will go through the page_* functions. In
39    // particular, we can now map memory using boot_map_region
40    // or page_insert
41    page_init();
42
43    check_page_free_list(1);
```

```

44     check_page_alloc();
45     check_page();
46
47     //////////////////////////////////////
48     // Now we set up virtual memory
49
50     //////////////////////////////////////
51     // Map 'pages' read-only by the user at linear address UPAGES
52     // Permissions:
53     //   - the new image at UPAGES -- kernel R, user R
54     //   (ie. perm = PTE_U | PTE_P)
55     //   - pages itself -- kernel RW, user NONE
56     boot_map_region(kern_pgdir, UPAGES, PTSIZE, PADDR(pages), PTE_P | PTE_W);
57
58     //////////////////////////////////////
59     // Use the physical memory that 'bootstack' refers to as the kernel
60     // stack. The kernel stack grows down from virtual address KSTACKTOP.
61     // We consider the entire range from [KSTACKTOP-PTSIZE, KSTACKTOP)
62     // to be the kernel stack, but break this into two pieces:
63     //   * [KSTACKTOP-KSTKSIZE, KSTACKTOP) -- backed by physical memory
64     //   * [KSTACKTOP-PTSIZE, KSTACKTOP-KSTKSIZE) -- not backed; so if
65     //     the kernel overflows its stack, it will fault rather than
66     //     overwrite memory. Known as a "guard page".
67     // Permissions: kernel RW, user NONE
68     boot_map_region(kern_pgdir, KSTACKTOP - KSTKSIZE, KSTKSIZE, PADDR(bootstack), PTE_P | PTE_W);
69
70     //////////////////////////////////////
71     // Map all of physical memory at KERNBASE.
72     // Ie. the VA range [KERNBASE, 2^32) should map to
73     // the PA range [0, 2^32 - KERNBASE)
74     // We might not have 2^32 - KERNBASE bytes of physical memory, but
75     // we just set up the mapping anyway.
76     // Permissions: kernel RW, user NONE
77     boot_map_region(kern_pgdir, KERNBASE, -KERNBASE, 0, PTE_P | PTE_W);
78
79
80     // Check that the initial page directory has been set up correctly.
81     check_kern_pgdir();
82
83     uint32_t cr4 = rcr4();
84     cr4 |= CR4_PSE;
85     lcr4(cr4);
86
87     // Switch from the minimal entry page directory to the full kern_pgdir
88     // page table we just created. Our instruction pointer should be
89     // somewhere between KERNBASE and KERNBASE+4MB right now, which is
90     // mapped the same way by both page tables.
91     //
92     // If the machine reboots at this point, you've probably set up your
93     // kern_pgdir wrong.
94     lcr3(PADDR(kern_pgdir));
95
96     check_page_free_list(0);
97
98     // entry.S set the really important flags in cr0 (including enabling
99     // paging). Here we configure the rest of the flags that we care about.
100     cr0 = rcr0();
101     cr0 |= CR0_PE|CR0_PG|CR0_AM|CR0_WP|CR0_NE|CR0_MP;
102     cr0 &= ~(CR0_TS|CR0_EM);
103     lcr0(cr0);
104
105     // Some more checks, only possible after kern_pgdir is installed.
106     check_page_installed_pgdir();
107 }

```

Question

Question 2

What entries (rows) in the page directory have been filled in at this point? What addresses do they map and where do they point? In other words, fill out this table as much as possible:

idx	va	pa	comment
957	ef400000	f011b000	pgdir self loop
956	ef000000	11c000	maps to UPAGES
959	efff8000	10f000	maps to bootstack
960	f0000000	0	maps to physical memory

idx	va	pa	comment
961	f0400000	400000	maps to physical memory
962	f0800000	800000	maps to physical memory
963	f0c00000	c00000	maps to physical memory
964	f1000000	1000000	maps to physical memory
965	f1400000	1400000	maps to physical memory
966	f1800000	1800000	maps to physical memory
967	f1c00000	1c00000	maps to physical memory
968	f2000000	2000000	maps to physical memory
969	f2400000	2400000	maps to physical memory
970	f2800000	2800000	maps to physical memory
971	f2c00000	2c00000	maps to physical memory
972	f3000000	3000000	maps to physical memory
973	f3400000	3400000	maps to physical memory
974	f3800000	3800000	maps to physical memory
975	f3c00000	3c00000	maps to physical memory
976	f4000000	4000000	maps to physical memory
977	f4400000	4400000	maps to physical memory
978	f4800000	4800000	maps to physical memory
979	f4c00000	4c00000	maps to physical memory
980	f5000000	5000000	maps to physical memory
981	f5400000	5400000	maps to physical memory
982	f5800000	5800000	maps to physical memory
983	f5c00000	5c00000	maps to physical memory
984	f6000000	6000000	maps to physical memory
985	f6400000	6400000	maps to physical memory
986	f6800000	6800000	maps to physical memory
987	f6c00000	6c00000	maps to physical memory
988	f7000000	7000000	maps to physical memory
989	f7400000	7400000	maps to physical memory
990	f7800000	7800000	maps to physical memory
991	f7c00000	7c00000	maps to physical memory

idx	va	pa	comment
992	f8000000	8000000	maps to physical memory
993	f8400000	8400000	maps to physical memory
994	f8800000	8800000	maps to physical memory
995	f8c00000	8c00000	maps to physical memory
996	f9000000	9000000	maps to physical memory
997	f9400000	9400000	maps to physical memory
998	f9800000	9800000	maps to physical memory
999	f9c00000	9c00000	maps to physical memory
1000	fa000000	a000000	maps to physical memory
1001	fa400000	a400000	maps to physical memory
1002	fa800000	a800000	maps to physical memory
1003	fac00000	ac00000	maps to physical memory
1004	fb000000	b000000	maps to physical memory
1005	fb400000	b400000	maps to physical memory
1006	fb800000	b800000	maps to physical memory
1007	fb000000	bc00000	maps to physical memory
1008	fc000000	c000000	maps to physical memory
1009	fc400000	c400000	maps to physical memory
1010	fc800000	c800000	maps to physical memory
1011	fcc00000	cc00000	maps to physical memory
1012	fd000000	d000000	maps to physical memory
1013	fd400000	d400000	maps to physical memory
1014	fd800000	d800000	maps to physical memory
1015	fdc00000	dc00000	maps to physical memory
1016	fe000000	e000000	maps to physical memory
1017	fe400000	e400000	maps to physical memory
1018	fe800000	e800000	maps to physical memory
1019	fec00000	ec00000	maps to physical memory
1020	ff000000	f000000	maps to physical memory
1021	ff400000	f400000	maps to physical memory
1022	ff800000	f800000	maps to physical memory

idx	va	pa	comment
1023	ffc00000	fc00000	maps to physical memory

Question 3

We have placed the kernel and user environment in the same address space. Why will user programs not be able to read or write the kernel's memory?

What specific mechanisms protect the kernel memory? We map kernel memory but not set the `PTE_U` flag.

Question 4

What is the maximum amount of physical memory that this operating system can support? Why?

4G, because the memory space is 32bit unsigned integer. The max value of a 32bit unsigned integer is $2^{32} = 4\text{G}$.

Question 5

How much space overhead is there for managing memory, if we actually had the maximum amount of physical memory? How is this overhead broken down? The overhead includes 2^{10} pagedir and 2^{10} page tables. They cost $(2^{10}+1) \cdot \text{PAGE_SIZE} = 2^{22} + 2^{12} = 4\text{M} + 4\text{k}$.

Question 6

At what point do we transition to running at an EIP above KERNBASE?

At `entry.S`, line 64

```
1  mov $relocated, %eax
2  jmp *%eax
3  relocated:
```

What makes it possible for us to continue executing at a low EIP between when we enable paging and when we begin running at an EIP above KERNBASE?

At `entrypgdir.c`, line 21

```
1  pde_t entry_pgdir[NPDENTRIES] = {
2      // Map VA's [0, 4MB) to PA's [0, 4MB)
3      [0]
4          = ((uintptr_t)entry_pgtbl - KERNBASE) + PTE_P,
5      // Map VA's [KERNBASE, KERNBASE+4MB) to PA's [0, 4MB)
6      [KERNBASE >> PDXSHIFT]
7          = ((uintptr_t)entry_pgtbl - KERNBASE) + PTE_P + PTE_W
8  };
```

The `entry_pgdir` maps VA's [0, 4MB) to PA's [0, 4MB), so code fetch won't crash.

Why is this transition necessary?

Because the map from VA's [0, 4MB) to PA's [0, 4MB) is a temporary map. When the virtual memory setups competely, this map won't exists. So the transition is necessary.

Challenge 1

We should turn on `CR4_PSE` to enable big page mode.

```
1 uint32_t cr4 = rcr4();
2 cr4 |= CR4_PSE;
3 lcr4(cr4);
```

My page mapping strategy is, in `boot_map_region`, big page is preferred.

```
1 static void
2 boot_map_region(pde_t *pgdir, uintptr_t va, size_t size, physaddr_t pa, int perm)
3 {
4     perm = perm | PTE_P;
5     size_t offset = 0;
6     while(offset < size) {
7         void * map_va = (void *)va + offset;
8         physaddr_t map_pa = pa + offset;
9         if(PTX(map_pa) == 0 && size - offset >= PGSIZE * NPENTRIES) {
10             pde_t * ppde = pgdir_walk_bigpg(pgdir, map_va);
11             if(ppde != NULL) {
12                 *ppde = map_pa | perm | PTE_PS;
13                 offset += PGSIZE * NPENTRIES;
14                 continue;
15             }
16         }
17         if(size - offset >= PGSIZE) {
18             pte_t * ppte = pgdir_walk(pgdir, map_va, true);
19             if(ppte == NULL) panic("No Available Page");
20             *ppte = map_pa | perm;
21             offset += PGSIZE;
22         }
23     }
24 }
```

Then, in `pgdir_walk`, big pages may be split as required.

```
1 pte_t *
2 pgdir_walk(pde_t *pgdir, const void *va, int create) {
3     pde_t * pde = &pgdir[PDX(va)];
4     if(*pde & PTE_P) {
5         if(*pde & PTE_PS) {
6             pte_t * ptab = split_large_page(pde);
7             if(ptab == NULL) panic("No Available Page");
8             return &ptab[PTX(va)];
9         }
10        else {
11            pte_t * ptab = KADDR(PDE_ADDR(*pde));
12            return &ptab[PTX(va)];
13        }
14    }
15    else if(create) {
16        struct PageInfo * ptab_info = page_alloc(ALLOC_ZERO);
17        if(ptab_info == NULL) return NULL;
18        ptab_info->pp_ref++;
19        physaddr_t pa = page2pa(ptab_info);
20        *pde = pa | PTE_P | PTE_U | PTE_W;
21        pte_t * ptab = KADDR(pa);
22        return &ptab[PTX(va)];
23    }
24    else return NULL;
25 }
```

The code for splitting page is:

```
1 pte_t * split_large_page(pde_t * pde) {
2     struct PageInfo * ptab_info = page_alloc(ALLOC_ZERO);
3     if(ptab_info == NULL) return NULL;
4     uint32_t flags = PDE_FLAGS(*pde) & ~PTE_PS;
5     physaddr_t pa = PDE_ADDR(*pde);
6     ptab_info->pp_ref++;
7     physaddr_t ptab_pa = page2pa(ptab_info);
8     pte_t * ptab = KADDR(ptab_pa);
9     for(size_t i = 0; i < NPENTRIES; i++) {
10         ptab[i] = (pa + i * PGSIZE) | flags;
```

```

11     }
12     *pde = ptab_pa | PTE_P | PTE_U | PTE_W;
13     return ptab;
14 }

```

The `check_va2pa` function is wrong when in big page mode. So I modified it:

```

1 static physaddr_t
2 check_va2pa(pde_t *pgdir, uintptr_t va)
3 {
4     pte_t *p;
5
6     pgdir = &pgdir[PDX(va)];
7     if (!(*pgdir & PTE_P))
8         return ~0;
9     if (*pgdir & PTE_PS) {
10        return PDE_ADDR(*pgdir) + PTX(va) * PGSIZE;
11    }
12    p = (pte_t*) KADDR(PTE_ADDR(*pgdir));
13    if (!p[PTX(va)] & PTE_P)
14        return ~0;
15    return PTE_ADDR(p[PTX(va)]);
16 }

```

Challenge 2

I implemented the follow commands:

```

1 mem pde # show all pde
2 mem show 0x0000 0xf000 # show pages from 0x0000 to 0xf000
3 mem set pde 0x0000 PS 1 # set PS flag of pde at 0x0000
4 mem dump 0xf0000000 0xf000f000 # dump virtual memory
5 mem dumpphy 0x0000 0xf000 # dump physical memory

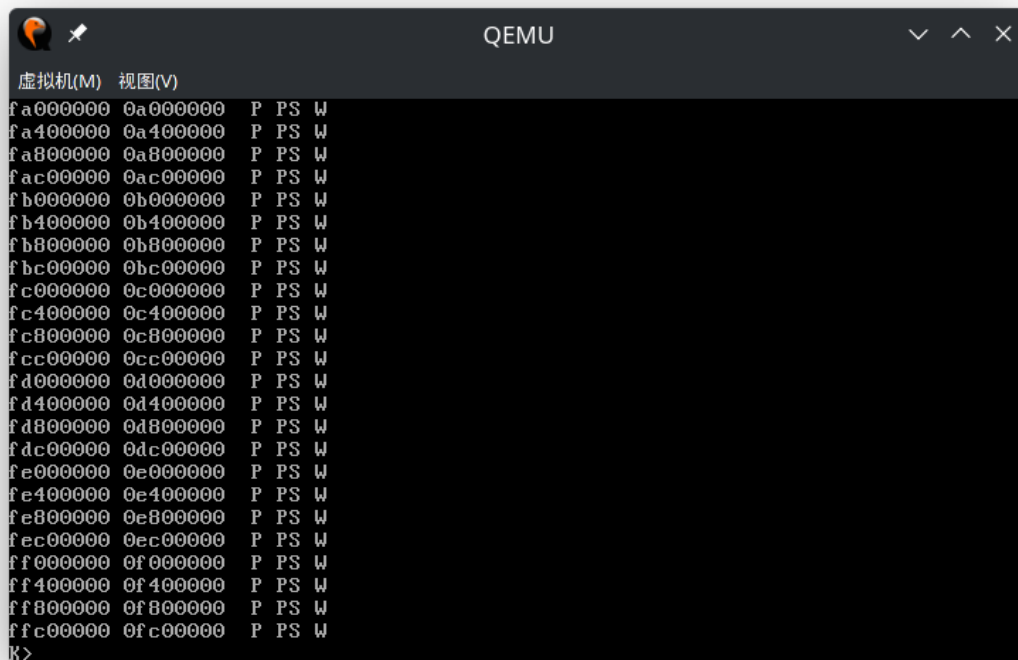
```

The command result:

```

1 mem pde

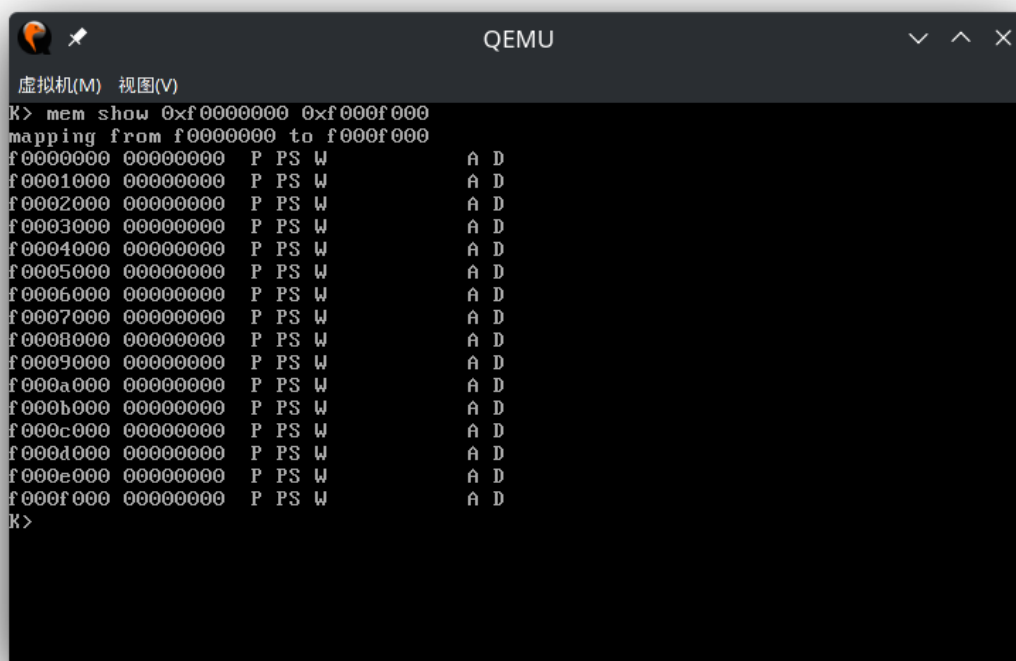
```



```

1 mem show 0xf0000000 0xf000f000

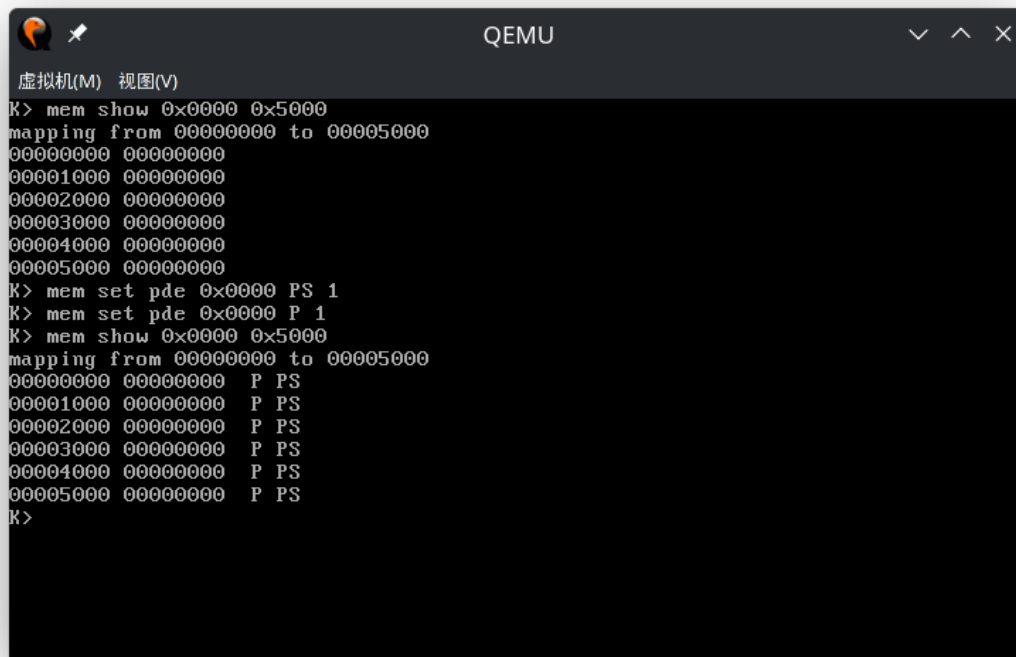
```



A screenshot of a QEMU terminal window. The title bar says 'QEMU'. Inside the terminal, the user has entered the command 'K> mem show 0xf0000000 0xf000f000'. The output shows a mapping from f0000000 to f000f000. The mapping is displayed as a table with columns for address, physical address, permissions, and device type. The addresses range from f0000000 to f000f000 in increments of 1000. The physical addresses are all 00000000. The permissions are all 'P PS W'. The device type is 'A D'.

```
K> mem show 0xf0000000 0xf000f000
mapping from f0000000 to f000f000
f0000000 00000000 P PS W      A D
f0001000 00000000 P PS W      A D
f0002000 00000000 P PS W      A D
f0003000 00000000 P PS W      A D
f0004000 00000000 P PS W      A D
f0005000 00000000 P PS W      A D
f0006000 00000000 P PS W      A D
f0007000 00000000 P PS W      A D
f0008000 00000000 P PS W      A D
f0009000 00000000 P PS W      A D
f000a000 00000000 P PS W      A D
f000b000 00000000 P PS W      A D
f000c000 00000000 P PS W      A D
f000d000 00000000 P PS W      A D
f000e000 00000000 P PS W      A D
f000f000 00000000 P PS W      A D
K>
```

```
1 mem show 0x0000 0x5000
2 mem set pde 0x0000 PS 1
3 mem set pde 0x0000 P 1
4 mem show 0x0000 0x5000
```



A screenshot of a QEMU terminal window. The title bar says 'QEMU'. Inside the terminal, the user has entered the command 'K> mem show 0x0000 0x5000'. The output shows a mapping from 00000000 to 00005000. The mapping is displayed as a table with columns for address, physical address, permissions, and device type. The addresses range from 00000000 to 00005000 in increments of 1000. The physical addresses are all 00000000. The permissions are all 'P PS'. The device type is 'A D'.

```
K> mem show 0x0000 0x5000
mapping from 00000000 to 00005000
00000000 00000000
00001000 00000000
00002000 00000000
00003000 00000000
00004000 00000000
00005000 00000000
K> mem set pde 0x0000 PS 1
K> mem set pde 0x0000 P 1
K> mem show 0x0000 0x5000
mapping from 00000000 to 00005000
00000000 00000000 P PS
00001000 00000000 P PS
00002000 00000000 P PS
00003000 00000000 P PS
00004000 00000000 P PS
00005000 00000000 P PS
K>
```

- Here, I set PS flag to 1, the pages from 0x0000 to 0x5000 is in the same big page.

```
1 mem dump 0xf0000000 0xf0000020
2 mem dumpphy 0x0000 0x0020
```

The screenshot shows a QEMU terminal window with the following output:

```
K> mem dump 0xf0000000 0xf0000020
53 ff 00 f0 53 ff 00 f0 c3 e2 00 f0 53 ff 00 f0
53 ff 00 f0 54 ff 00 f0 53 ff 00 f0 53 ff 00 f0
a5
K> mem dumpphy 0x0000 0x0020
53 ff 00 f0 53 ff 00 f0 c3 e2 00 f0 53 ff 00 f0
53 ff 00 f0 54 ff 00 f0 53 ff 00 f0 53 ff 00 f0
a5
K>
```

The code is:

```
1 // -----
2 // Debug functions.
3 // -----
4
5 static int64_t atoi(char * arg) {
6 #define _isnum(x) (('0' <= (x) && (x) <= '9') || ('a' <= (x) && (x) <= 'f')
7 #define _getnum(x) (('0' <= (x) && (x) <= '9') ? ((x) - '0') : ((x) - 'a' + 10))
8 int64_t ret = 0, mul = 1, base=10;
9 while(!('0' <= *arg && *arg <= '9')) arg++;
10 if(*arg == '-') { mul = -1; arg++; }
11 if(arg[0] == '0') {
12     if(arg[1] == 'x') {
13         base = 16;
14         arg = arg + 2;
15     }
16 }
17 while(_isnum(*arg)) {
18     ret = ret * base + _getnum(*arg);
19     arg++;
20 }
21 return ret * mul;
22 #undef _getnum
23 #undef _isnum
24 }
25
26 static void _show_pte(pte_t pde) {
27 #define _show_pde_inner(flag) \
28     if(pde & PTE_##flag) {cprintf("#flag"); cprintf(" ");} \
29     else {cprintf(" ");} for(size_t i = strlen("#flag"); i-->0; i--) cprintf(" ");
30 cprintf("%08x ", PTE_ADDR(pde));
31 _show_pde_inner(P);
32 _show_pde_inner(PS);
33 _show_pde_inner(W);
34 _show_pde_inner(U);
35 _show_pde_inner(PWT);
36 _show_pde_inner(PCD);
37 _show_pde_inner(A);
38 _show_pde_inner(D);
39 _show_pde_inner(G);
40 _show_pde_inner(AVAIL);
41 #undef _show_pde_inner
42 }
43
```

```

44 static void _show_pde(pde_t pde) {
45     _show_pte((pte_t)pde);
46 }
47
48 int
49 memcmd_pde(int argc, char ** argv, struct Trapframe * tf) {
50     for(size_t i = 0; i < NPENTRIES; i++) {
51         uintptr_t va = PGSIZE * NPENTRIES * i;
52         cprintf("%08x ", va);
53         _show_pde(kern_pgdir[PDX(va)]);
54         cprintf("\n");
55     }
56     return 0;
57 }
58
59 int
60 memcmd_show(int argc, char ** argv, struct Trapframe * tf) {
61     if(argc != 3) {
62         cprintf("usage mem show <start> <end>\n");
63         return -1;
64     }
65     uintptr_t start = ROUNDDOWN(atoi(argv[1]), PGSIZE);
66     uintptr_t end = ROUNDUP(atoi(argv[2]) + 1, PGSIZE);
67     cprintf("mapping from %08x to %08x\n", start, end - PGSIZE);
68     uintptr_t va = start;
69     do {
70         pde_t * ppde = &kern_pgdir[PDX(va)];
71         if(!(*ppde & PTE_P) || (*ppde & PTE_PS)) {
72             cprintf("%08x ", va); _show_pde(*ppde); cprintf("\n");
73         }
74         else {
75             pte_t * ppte = pgdir_walk(kern_pgdir, (void*)va, false);
76             assert(ppte);
77             cprintf("%08x ", va); _show_pte(*ppte); cprintf("\n");
78         }
79         va += PGSIZE;
80     } while(va != end);
81     return 0;
82 }
83
84 int
85 memcmd_set(int argc, char ** argv, struct Trapframe * tf) {
86     if(argc != 5) {
87         cprintf("usage mem set <pde|pte> <va> <ent> <value>\n");
88         return -1;
89     }
90     bool pde = strcmp(argv[1], "pde") == 0;
91     uintptr_t va = atoi(argv[2]);
92     const char * ent = argv[3];
93     uint64_t value = atoi(argv[4]);
94     #define _memcmd_set_inner(flag) \
95     else if(strcmp(ent, #flag)==0) {\
96         if(pde) {\
97             pde_t * ppde = &kern_pgdir[PDX(va)]; \
98             *ppde = (*ppde & (~PTE_##flag)) | (value * PTE_##flag);\
99         } \
100     else {\
101         pte_t * ppte = pgdir_walk(kern_pgdir, (void*)va, true); \
102         if(ppte == NULL) {cprintf("No Available Pages"); \
103             return -1; \
104         } \
105         *ppte = (*ppte & (~PTE_##flag)) | (value * PTE_##flag); \
106     } \
107 }
108 if(strcmp(ent, "pa")==0) {
109     if(pde) {
110         pde_t * ppde = &kern_pgdir[PDX(va)];
111         *ppde = value | PDE_FLAGS(*ppde);
112         tlb_invalidate(kern_pgdir, (void*)va);
113     }
114     else {
115         pte_t * ppte = pgdir_walk(kern_pgdir, (void*)va, true);
116         *ppte = value | PTE_FLAGS(*ppte);
117         tlb_invalidate(kern_pgdir, (void*)va);
118     }
119 }
120 _memcmd_set_inner(P)
121 _memcmd_set_inner(PS)
122 _memcmd_set_inner(W)
123 _memcmd_set_inner(U)
124 _memcmd_set_inner(PWT)
125 _memcmd_set_inner(PCD)
126 _memcmd_set_inner(A)
127 _memcmd_set_inner(D)
128 _memcmd_set_inner(G)
129 _memcmd_set_inner(AVAIL)
130 else {
131     cprintf("Unknown entity: %s\n", ent);
132     return -1;
133 }
134 return 0;
135 }
136
137 int
138 memcmd_dump(int argc, char ** argv, struct Trapframe * tf) {
139     if(argc != 3) {
140         cprintf("usage mem dump <start> <end>\n");
141         return -1;
142     }

```



```
143     uintptr_t start = atoi(argv[1]);
144     uintptr_t end = atoi(argv[2]);
145     short tick = 0;
146     for(uintptr_t i = start; i <= end; i++) {
147         printf("%02x ", *(unsigned char *)i);
148         if(((++tick)&0xf)==0) printf("\n");
149     }
150     if(tick)printf("\n");
151     return 0;
152 }
153
154 int
155 memcmd_dumpphy(int argc, char ** argv, struct Trapframe * tf) {
156     if(argc != 3) {
157         printf("usage mem dumpphy <start> <end>\n");
158         return -1;
159     }
160     uintptr_t start = atoi(argv[1]);
161     uintptr_t end = atoi(argv[2]);
162     short tick = 0;
163     for(uintptr_t i = start; i <= end; i++) {
164         printf("%02x ", *(unsigned char*)KADDR(i));
165         if(((++tick)&0xf)==0) printf("\n");
166     }
167     if(tick)printf("\n");
168     return 0;
169 }
170
171 int
172 mem_memcmd(int argc, char ** argv, struct Trapframe * tf) {
173     if(argc == 1) {
174         printf("Usage: mem <pde|show|set|dump|dumpphy> ...");
175         return -1;
176     }
177     else {
178         if(strcmp(argv[1], "pde") == 0) {
179             memcmd_pde(argc - 1, argv + 1, tf);
180         }
181         else if(strcmp(argv[1], "show") == 0) {
182             memcmd_show(argc - 1, argv + 1, tf);
183         }
184         else if(strcmp(argv[1], "set") == 0) {
185             memcmd_set(argc - 1, argv + 1, tf);
186         }
187         else if(strcmp(argv[1], "dump") == 0) {
188             memcmd_dump(argc - 1, argv + 1, tf);
189         }
190         else if(strcmp(argv[1], "dumpphy") == 0) {
191             memcmd_dumpphy(argc - 1, argv + 1, tf);
192         }
193         else {
194             printf("Unknown command %s\n", argv[1]);
195             return -1;
196         }
197     }
198     return 0;
199 }
```