
Lab #2

陈威宇

Exercise 1

boot_alloc 分配 n 字节物理内存，只在 JOS 创建虚拟内存系统的时候调用，之后就不用了。代码如下，注意判断 out of memory.

```
1  if (n>0){
2      uint32_t sz = ROUNDUP(n, PGSIZE);
3      nextfree += sz;
4      if (((uintptr_t)nextfree)-KERNBASE > 0x00400000)
5          panic("out of memory");
6      return (nextfree-sz);
7  }
8  else{
9      return nextfree;
10 }
```

mem_init 里的 npages 是物理页的总数, pages 是维护这些物理页的数组, 每个元素的类型都是 struct PageInfo.

```
1  pages=(struct PageInfo *)boot_alloc(npages*sizeof(struct PageInfo));
2  memset(pages,0,npages*sizeof(struct PageInfo));
```

page_init 的作用是把 pages 数组初始化, 并初始化维护 free physical page 的 page_free_list 链表.

```
1  void
2  page_init(void)
3  {
4      page_free_list=NULL;
5      size_t i;
6      for (i = 0; i < npages; i++) {
7          pages[i].pp_ref = 0;
8          if (i==0 || (i*PGSIZE >= IOPHYSMEM && i*PGSIZE < EXTPHYSMEM) || (i*PGSIZE>=0x100000 && i*P
9              pages[i].pp_link = NULL;
10     }
11     else{
12         pages[i].pp_link = page_free_list;
13         page_free_list = &pages[i];
14     }
```

```
15     }
16 }
```

page_alloc 的作用是分配一个 free physical page, 直接从 free list 取即可.

```
1  struct PageInfo *
2  page_alloc(int alloc_flags)
3  {
4      // Fill this function in
5      struct PageInfo * t;
6      if (page_free_list==NULL) return NULL;
7      t = page_free_list;
8      page_free_list = (t -> pp_link);
9      if ((t->pp_ref)!=0){
10         panic("PANIC!!!");
11     }
12     (t->pp_link)=NULL;
13     if (alloc_flags & ALLOC_ZERO)
14         memset(page2kva(t),0,PGSIZE);
15     return t;
16 }
```

page_free 的用处是将一个 page 放回 free list, 当然, 要确保其被 reference 的次数已经是 0 了.

```
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)!=0 || (pp->pp_link)!=NULL){
8         panic("page_free PANIC!!!");
9     }
10     (pp->pp_link) = page_free_list;
11     page_free_list = pp;
12 }
```

Question after Exercise 3

1

变量 x 的类型应当是 uintptr_t.

因为 *value 可以直接修改, 所以 value 是虚拟地址.

Exercise 4

`pgdir_walk` 的作用是定位到一个 virtual address 对应的 PTE，直接查 page directory 即可。如果其对应的 page table 不存在的话，根据 `create` 的值可能会新建一个。

```
1 pte_t *
2 pgdir_walk(pde_t *pgdir, const void *va, int create)
3 {
4     // Fill this function in
5     uintptr_t la = (uintptr_t)va;
6     pde_t t = (*(pgdir+PDX(la)));
7     if (t&PTE_P){
8         physaddr_t p1 = PTE_ADDR(t);
9         uintptr_t p2 = (uintptr_t)KADDR(p1);
10        return ((pte_t *)p2)+PTX(la);
11    }
12    if (create==0)
13        return NULL;
14    struct PageInfo * tmp = page_alloc(1);
15    if (tmp==NULL)
16        return NULL;
17    memset(page2kva(tmp),0,PGSIZE);
18    ++(tmp->pp_ref);
19    (*(pgdir+PDX(la))) = (page2pa(tmp) | 0x007);
20    return ((pte_t *)page2kva(tmp)) + PTX(la) ;
21 }
```

`boot_map_region` 的作用是，将连续的一些 virtual page 映射到连续的一些 physical page。做法是，对每个 virtual page，通过 `pgdir_walk` 找到 PTE，然后直接修改即可。

```
1 static void
2 boot_map_region(pde_t *pgdir, uintptr_t va, size_t size, physaddr_t pa, int perm)
3 {
4     // Fill this function in
5     if (size==0) return;
6     for (int i=0;i<size/PGSIZE;++i){
7         uintptr_t VA=va+i*PGSIZE;
8         physaddr_t PA=pa+(VA-va);
9         pte_t * tmp = pgdir_walk(pgdir, (void *)VA, 1);
10        if (tmp==NULL)
11            panic("PANIC!!!");
12        (*tmp) = PA | (perm|PTE_P);
13    }
```

```
13     }
14 }
```

page_lookup 的作用是对一个 virtual address 求其对应的物理页. 调用 pgdir_walk 后直接根据 PTE 上的 physical address 转换为物理页即可 (利用 pa2page).

```
1  struct PageInfo *
2  page_lookup(pde_t *pgdir, void *va, pte_t **pte_store)
3  {
4      // Fill this function in
5      pte_t * t = pgdir_walk(pgdir, va, 0);
6      if (t==NULL)
7          return NULL;
8      if (!((*t)&PTE_P))
9          return NULL;
10     if (pte_store!=0){
11         (*pte_store) = t;
12     }
13     return pa2page(PTE_ADDR(*t));
14 }
```

page_remove 的作用是, 将 va 这个 virtual address 对应的 virtual page 取消映射. 直接修改 PTE 即可, 注意还要将映射到的 physical page 的 ref 减 1, 即 decref. 另外还要 invalidate TLB.

```
1  void
2  page_remove(pde_t *pgdir, void *va)
3  {
4      // Fill this function in
5      pte_t * tmp;
6      struct PageInfo * pg = page_lookup(pgdir, va, &tmp);
7      if (pg==NULL)
8          return;
9      page_decref(pg);
10     if (tmp!=NULL)
11         (*tmp)=0;
12     tlb_invalidate(pgdir, va);
13 }
```

page_insert 的作用是, 将一个 virtual address 对应的 virtual page 映射到物理页 pp. 如果原先就有映射, 要先 remove 掉. 需要注意的一点是, 如果原先映射到的就是 pp, 不要在 remove 的时候将其加到 free list 上, 我处理这个细节的办法是, 在调用 page_remove 前给 pp->ref 加 1, 调用后再减 1.

```
1  int
2  page_insert(pde_t *pgdir, struct PageInfo *pp, void *va, int perm)
```

```

3 {
4     // Fill this function in
5     pte_t * t = pgdir_walk(pgdir, va, 1);
6     if (t==NULL)
7         return -E_NO_MEM;
8     if ((*t)&PTE_P){
9         ++(pp->pp_ref);
10        page_remove(pgdir, va);
11        --(pp->pp_ref);
12    }
13    (*t) = page2pa(pp) | (perm|PTE_P);
14    ++(pp->pp_ref);
15    return 0;
16 }

```

Exercise 5

直接调用 `boot_map_region` 来建立内存映射.

```

1 boot_map_region(kern_pgdir, UPAGES, PTSIZE, PADDR(pages), PTE_U | PTE_P);
2 boot_map_region(kern_pgdir, KSTACKTOP-KSTKSIZE, KSTACKTOP-(KSTACKTOP-KSTKSIZE), PADDR(bootst
3 boot_map_region(kern_pgdir, KERNBASE, 0xffffffff-KERNBASE+1, 0, PTE_W | PTE_P);

```

Questions after Exercise 5

2

Entry	Base Virtual Address	Points to (logically):
-----	-----	-----
1023	0xffc00000	Page table for top 4MB of phys memory
.....
960	0xf0000000	Page table for bottom 4MB of phys memory
959	0xefc00000	Page table for kernel stack
958	0xef800000	NULL
957	0xef400000	Page table for Page Directory
956	0xef000000	Page table for RO PAGES
955	0xeec00000	NULL
.....
1	0x00400000	NULL
0	0x00000000	NULL

3

对于 kernel momoey, 将 Page table/directory entry flags 的 PTE_U 位设为 0, 用户就不能访问内核内存了。

4

256M. 这是因为, JOS 把所有物理内存都 map 到了虚拟地址 KERNBASE. 而 KERNBASE=0xf0000000, 于是物理内存最多为 0x100000000 字节, 即 256M。

5

此时 overhead 分为 3 部分.

1. 存放 page directory 的空间, 即 4KB. 2. 存放所有 page table 的空间, 即 $1024 * 4KB = 4MB$ 3. 维护 physical page 的数组 pages, 占用空间 $\text{sizeof}(\text{struct PageInfo}) * \text{npages}$, 即 $8B * (256MB/4K) = 0.5MB$, 所以总的 overhead 有 $4.5MB + 4KB$.

6

执行完 entry.S 里的 `jmp *%eax` 语句后 EIP 跳转到比 KERNBASE 大的地方. 在 turn on paging 与 `jmp *%eax` 之间, 能正常执行是因为 `entry_pgdir` 建立了虚拟地址 $[0, 4M)$ 到物理地址 $[0, 4M)$ 的映射. 这个转变是必要的, 因为 kernel 要在虚拟地址的高地址执行, 把低地址留给用户, 起到对 kernel memory 的保护作用.