

练习0：填写已有实验

将 lab1 的 `kern/debug/kdebug.c`、`kern/init/init.c` 以及 `kern/trap/trap.c` 直接复制到 lab2 里即可

练习1：实现 first-fit 连续物理内存分配算法

因为完全不会，所以去 csdn 找了几篇文章先看了看

list_entry 结构体

写于：libs/list.h

```
1 struct list_entry {
2     struct list_entry *prev, *next;
3 };
4
5 typedef struct list_entry list_entry_t;
```

这里定义了一个双向链表

Page 结构体

写于：kern/mm/memlayout.h

```
1 struct Page {
2     int ref;           // page frame's reference counter
3     uint32_t flags;    // array of flags that describe the status of the
4     // page frame
5     unsigned int property; // the num of free block, used in first fit pm
6     // manager
7     list_entry_t page_link; // free list link
8 };
```

1. ref 表示的是：这个页被页表的引用记数，也就是映射此物理页的虚拟页个数

如果这个页被页表引用了，即在某页表中有一个页表项设置了一个虚拟页到这个 Page 管理的物理页的映射关系，就会把 Page 的 ref 加一

反之，若页表项取消，即映射关系解除，就会把 Page 的 ref 减一

2. flags 表示此物理页的状态标记，有两个标志位状态

为 1 的时候，代表这一页是 free 状态，可以被分配，但不能对它进行释放

如果为 0，那么说明这个页已经分配了，不能被分配，但是可以被释放掉

3. property 用来记录某连续空闲页的数量

这里需要注意的是用到此成员变量的这个 Page **一定是连续内存块的开始地址（第一页的地址）**

4. page_link 是便于把多个连续内存空闲块链接在一起的双向链表指针

连续内存空闲块利用这个页的成员变量 page_link 来链接比它地址小和大的其他连续内存空闲块
释放的时候只要将这个空间通过指针放回到双向链表中

free_area_t 结构体

写于: kern/mm/mmlayout.h

```
1 typedef struct {
2     list_entry_t free_list; // the list header
3     unsigned int nr_free;   // # of free pages in this free list
4 } free_area_t;
```

在初始情况下,也许这个物理内存的空闲物理页都是连续的,这样就形成了一个大的连续内存空闲块但随着物理页的分配与释放,这个大的连续内存空闲块会分裂为一系列地址不连续的多个小连续内存空闲块

且每个连续内存空闲块内部的物理页是连续的,那么为了有效地管理这些小连续内存空闲块

所有的连续内存空闲块可用一个双向链表管理起来,便于分配和释放,为此定义了一个 free_area_t 数据结构

包含了一个 list_entry 结构的双向链表指针和记录当前空闲页的个数的无符号整型变量 nr_free

其中的链表指针指向了空闲的物理页

该数据结构中含有两个成员:

free_list: 一个list_entry结构的双向链表指针

nr_free: 记录当前空闲页的个数

其中,在 kern/mm/pmm.c 定义了一个该数据结构的实例,用于管理实际分配当中的空闲空间

ucore 利用一个物理内存管理类 pmm_manager 确定使用需要的分配算法: kern/mm/pmm.h

pmm_manager 结构体

```
1 struct pmm_manager {
2     const char *name; // 物理内存页管理器的名字
3     void (*init)(void); // 初始化内存管理器
4     void (*init_memmap)(struct Page *base, size_t n); // 初始化管理空闲内存页的
    数据结构
5     struct Page *(*alloc_pages)(size_t n); // 分配 n 个物理内存页
6     void (*free_pages)(struct Page *base, size_t n); // 释放 n 个物理内存页
7     size_t (*nr_free_pages)(void); // 返回当前剩余的空闲页数
8     void (*check)(void); // 用于检测分配/释放实现是否正确的辅助函数
9 };
```

总的来说,一个物理内存管理类包含这个类的名字和一些函数,类似:初始化、分配页、释放页等操作可以看做类似于 lab6 中的调度的五元组的绑定方法

同样,我们所实现的函数,也是通过名字绑定的,具体的绑定过程在 kern/mm/pmm.c:

```

1  const struct pmm_manager default_pmm_manager = {
2      .name = "default_pmm_manager",
3      .init = default_init,
4      .init_memmap = default_init_memmap,
5      .alloc_pages = default_alloc_pages,
6      .free_pages = default_free_pages,
7      .nr_free_pages = default_nr_free_pages,
8      .check = default_check,
9  };

```

接下来，我们来实现 First fit 算法的相关函数：

default_init, default_init_memmap, default_alloc_pages, default_free_pages

default_init 函数

```

1  free_area_t free_area;
2
3  #define free_list (free_area.free_list)
4  #define nr_free (free_area.nr_free)
5
6  static void
7  default_init(void) {
8      list_init(&free_list);
9      nr_free = 0;
10 }

```

调用库函数 list_init 初始化掉 free_area_t

即管理所有连续的空闲内存空间块的数据结构 free_area_t 的双向链表和空闲块数

1. list_init 函数

写于：libs/list.h

```

1  static inline void
2  list_init(list_entry_t *elm) {
3      elm->prev = elm->next = elm;
4  }

```

大致作用就是初始化双向循环链表

这里有两个宏定义，在 kern/mm/pmm.c：

```

1  #define free_list (free_area.free_list)
2  #define nr_free (free_area.nr_free)

```

从这里看出来，初始化的对象是 free_area_t

这是与具体物理内存分配算法无关的，因此直接使用默认的函数实现即可

default_init_memmap 函数

init_memmap 函数主要实现的是一个根据现有的内存情况构建空闲块列表的初始状态的功能

何时应该执行这个函数呢，因此我们需要研究一下这个函数是如何被调用的：

调用过程是：kern_init → pmm_init → page_init → init_memmap

1. kern_init 函数

写于 kern/init/init.c

这个函数是进入 ucore 操作系统之后第一个执行的函数，对内核进行初始化

其中调用了初始化物理内存的函数 pmm_init

2. pmm_init 函数

写于：kern/mm/pmm.c

这个函数主要是完成对于整个物理内存的初始化，页初始化只是其中的一部分

调用位置偏前，函数之后的部分可以不管，直接进入 page_init 函数

3. page_init 函数

写于：kern/mm/pmm.c

```
1  static void
2  page_init(void) {
3  struct e820map *memmap = (struct e820map *) (0x8000 + KERNBASE);
4      // 首先声明一个e820map类的对象memmap，与物理内存相关
5      // 在本实验中，我们获取内存信息的方式是通过e820中断（一种特殊的内核中断模式）
6      uint64_t maxpa = 0;
7
8      cprintf("e820map:\n");
9      int i;
10     for (i = 0; i < memmap->nr_map; i++) {
11         // 这里可以看做一个遍历，第一轮遍历是遍历物理地址空间
12         // 获取物理地址的最大值maxpa（探测物理内存布局）
13         uint64_t begin = memmap->map[i].addr, end = begin + memmap->
14         map[i].size;
15         cprintf("  memory: %08llx, [%08llx, %08llx], type = %d.\n",
16                 memmap->map[i].size, begin, end - 1, memmap->
17         map[i].type);
18         if (memmap->map[i].type == E820_ARM) {
19             if (maxpa < end && begin < KMEMSIZE) {
20                 maxpa = end;
21             }
22         }
23         // maxpa不能大于所允许的最大值，这个最大值宏定义在memlayout的第57行
24         // 物理地址所允许的最大值为0x38000000 (KMEMSIZE)
25         if (maxpa > KMEMSIZE) {
26             maxpa = KMEMSIZE;
27         }
28
29         extern char end[];
30
31         // 创建的页数量等于物理地址最大值除以页大小，其中页大小为4096字节，即4KB
32         // 该定义在mmu.h的第226行
33         npage = maxpa / PGSIZE;
34         pages = (struct Page *) ROUNDUP((void *) end, PGSIZE);
35
36         // 将所有的页设置为保留页，在实际初始化页面init_memmap的时候，又会更改回非保留
37         // 推测是在初始化过程中这样处理，是为了防止页面被分配，结构被破坏。
38         for (i = 0; i < npage; i++) {
39             SetPageReserved(pages + i);
40         }
41     }
42 }
```

```

39     }
40
41     uintptr_t freemem = PADDR((uintptr_t)pages + sizeof(struct Page) *
npage);
42
43     // 第二次遍历物理内存，这一次遍历主要是调用init_memmap初始化各个页表
44     for (i = 0; i < memmap->nr_map; i++) {
45         uint64_t begin = memmap->map[i].addr, end = begin + memmap->
map[i].size;
46         if (memmap->map[i].type == E820_ARM) {
47             if (begin < freemem) {
48                 begin = freemem;
49             }
50             if (end > KMEMSIZE) {
51                 end = KMEMSIZE;
52             }
53             if (begin < end) {
54                 begin = ROUNDUP(begin, PGSIZE);
55                 end = ROUNDDOWN(end, PGSIZE);
56                 if (begin < end) {
57                     init_memmap(pa2page(begin), (end - begin) /
PGSIZE);
58                 }
59             }
60         }
61     }
62 }

```

page_init 函数主要是完成了一个**整体物理地址的初始化过程**，包括设置标记位，探测物理内存布局等操作

上面函数的注释中，标出了几个关键位置代码

但是，其中最关键的部分，也是和实验相关的页初始化，交给了 init_memmap 函数处理

这里我们主要研究一下上方 init_memmap 传入的两个参数：

1. pa2page(begin)

写于：kern/mm/pmm.h

```

1 static inline struct Page *
2 pa2page(uintptr_t pa) {
3     if (PPN(pa) >= npage) {
4         panic("pa2page called with invalid pa");
5     }
6     return &pages[PPN(pa)];
7 }

```

其中 PPN 是物理地址页号

该函数的作用是：返回传入参数 pa 开始的第一个物理页

其基地址是 base (default_init_memmap 函数里有写)

2. (end - begin) / PGSIZE

由于 end 和 begin 都是循环中记录位置的标记，PGSIZE 为 4 KB 的页大小，这里就代表物理页的个数

接下来看 default_init_memmap 函数：

```

1 static void
2 default_init_memmap(struct Page *base, size_t n) {
3     assert(n > 0);
4     struct Page *p = base;
5     for (; p != base + n; p++) {
6         assert(PageReserved(p));
7         p->flags = p->property = 0;
8         set_page_ref(p, 0);
9     }
10    base->property = n;
11    SetPageProperty(base);
12    nr_free += n;
13    list_add(&free_list, &(base->page_link));
14 }

```

这里先看几个不理解的函数：

1. set_page_ref 函数：

写于：kern/mm/pmm.h

```

1 static inline void
2 set_page_ref(struct Page *page, int val) {
3     page->ref = val;
4 }

```

这个函数就是设置记录映射此物理页的虚拟页的变量的数值

2. test_bit 函数：

写于：libs/atomic.h

```

1 static inline bool
2 test_bit(int nr, volatile void *addr) {
3     int oldbit;
4     asm volatile ("btl %2, %1; sbbl %0,%0" : "=r" (oldbit) : "m" (*
5     (volatile long *)addr), "Ir" (nr));
6     return oldbit != 0;

```

这里 bt 的作用是把 *addr 的第 nr 位复制到 cf

btl 里的 l 是用于操作 32 位的长字值所使用的后缀

sbb 与 sub 的区别：

sub ax, bx 的结果是 ax - bx

sbb ax, bx 的结果是 ax - bx - cf (进/借位标志)

这里 oldbit 初始化的时候不用管值是多少，最后经过 sbbl 汇编，只会有两个值

一个是 0，一个是 -1

那么最后检查 oldbit 的值是否为 0，若 oldbit 的值为 0，则 *addr 第 nr 位的值为 0，否则为 1

也就是说，这个函数返回 true 的话，该测试位上的值就为 0，反之则为 1

3. PageReserved 函数：

写于：kern/mm/memlayout.h

```
1 | #define PageReserved(page)      test_bit(PG_reserved, &((page)->flags))
```

接下来看 **PG_reserved** 的宏定义：

```
1 | #define PG_reserved              0          // if this bit=1: the Page
    | is reserved for kernel, cannot be used in alloc/free_pages; otherwise,
    | this bit=0
```

这个函数的意思就是检查 `(page)->flags` 这个值的第 0 位是否为 0

如果为 0，就可以顺利通过 `assert(PageReserved(p));` 语句

如果为 1，则报告错误并终止程序

Reserved 表示该页被操作系统保留，不能进行分配或者释放

那么该函数的意思就是如果该页不是被操作系统保留的，那么就不报错，顺利通过

4. **set_bit** 函数：

写于： **libs/atomic.h**

```
1 | /* *
2 |  * set_bit - Atomically set a bit in memory
3 |  * @nr:      the bit to set
4 |  * @addr:    the address to start counting from
5 |  *
6 |  * Note that @nr may be almost arbitrarily large; this function is not
7 |  * restricted to acting on a single-word quantity.
8 |  */
9 | static inline void
10 | set_bit(int nr, volatile void *addr) {
11 |     asm volatile ("btsl %1, %0" : "=m" (*(volatile long *)addr) : "Ir"
12 | (nr));
    | }
```

bts 在执行 bt 命令的同时（把 *addr 的第 nr 位复制到 cf），把操作数的指定位置设置为 1

这个函数的功能就是将 *addr 的第 nr 位设为 1

第一个占位符 %0 与 C 语言变量 addr 对应，第二个占位符 %1 与 C 语言变量 nr 对应

因此上面的汇编语句代码与此伪代码等价： `btsl nr, addr`

该指令的两个操作数不能全是内存变量，因此将 nr 的限定字符串指定为 **Ir**

将 nr 与立即数或者寄存器相关联，这样两个操作数中只有 addr 为内存变量

5. **SetPageProperty** 函数：

写于： **kern/mm/memlayout.h**

```
1 | #define SetPageProperty(page)      set_bit(PG_property, &((page)->flags))
```

接下来看 **PG_property** 的宏定义：

```

1  #define PG_property          1          // if this bit=1: the Page
    is the head page of a free memory block(contains some
    continuous_addrress pages), and can be used in alloc_pages; if this
    bit=0: if the Page is the the head page of a free memory block, then
    this Page and the memory block is allocated. Or this Page isn't the head
    page.

```

这个函数的意思就是修改 `(page)->flags` 这个值的第 1 位为 1

6. list_add 函数等:

写于: `libs/atomic.h`

```

1  static inline void
2  list_add(list_entry_t *listelm, list_entry_t *elm) {
3      list_add_after(listelm, elm);
4  }
5
6  static inline void
7  list_add_before(list_entry_t *listelm, list_entry_t *elm) {
8      __list_add(elm, listelm->prev, listelm);
9  }
10
11 static inline void
12 list_add_after(list_entry_t *listelm, list_entry_t *elm) {
13     __list_add(elm, listelm, listelm->next);
14 }
15
16 static inline void
17 __list_add(list_entry_t *elm, list_entry_t *prev, list_entry_t *next) {
18     prev->next = next->prev = elm;
19     elm->next = next;
20     elm->prev = prev;
21 }

```

所以能看出来 `list_add` 实质上就是 `list_add_after`, 它和 `list_add_before` 都调用了 `__list_add`

Answer:

```

1  static void
2  default_init_memmap(struct Page *base, size_t n) {
3      assert(n > 0);
4      struct Page *p = base;
5      for (; p != base + n; p++) {
6          assert(PageReserved(p));
7          p->flags = p->property = 0;
8          set_page_ref(p, 0);
9      }
10     base->property = n;
11     SetPageProperty(base);
12     nr_free += n;
13     list_add_before(&free_list, &(base->page_link));
14 }

```

这个函数的主要目的就是初始化所有空闲物理页的 Page 结构, 具体初始化方式为:

`assert(n > 0);` 确保不会有分配 0 页的情况发生

`struct Page *p = base;` 设置指针变量 `p`, 为了在之后遍历 `n` 页

`for (; p != base + n; p++)` 循环内容:

`assert(PageReserved(p));` 确保该页没有被操作系统保留 (检测 `p->flags` 的第 0 位)

`p->flags = p->property = 0;` 将 `Page` 结构中的两个成员变量的值置零, 可以当作手动初始化

将遍历到的物理页的状态标记设置为 0, 说明该页已经被分配了

`set_page_ref(p, 0);` 清空这些物理页的引用计数

`base->property = n;` 设置页头中记录连续空闲页数量的值为 `n`

`SetPageProperty(base);` 修改 `(base->flags)` 值的第 1 位为 1

`nr_free += n;` 更新总空闲块数目 (该变量一开始在 `default_init` 函数中初始化)

`list_add_before(&free_list, &(base->page_link));`

将 `&(base->page_link)` 插入到 `(&free_list->prev)` 和 `&free_list` 的中间

这里因为 `&free_list` 是双向链表, 所以插入是没有问题的, 用 `list_add` 也没问题

default_free_pages 函数

先看初始函数:

```
1  static void
2  default_free_pages(struct Page *base, size_t n) {
3      assert(n > 0);
4      struct Page *p = base;
5      for (; p != base + n; p++) {
6          assert(!PageReserved(p) && !PageProperty(p));
7          p->flags = 0;
8          set_page_ref(p, 0);
9      }
10     base->property = n;
11     SetPageProperty(base);
12     list_entry_t *le = list_next(&free_list);
13     while (le != &free_list) {
14         p = le2page(le, page_link);
15         le = list_next(le);
16         if (base + base->property == p) {
17             base->property += p->property;
18             ClearPageProperty(p);
19             list_del(&(p->page_link));
20         }
21         else if (p + p->property == base) {
22             p->property += base->property;
23             ClearPageProperty(base);
24             base = p;
25             list_del(&(p->page_link));
26         }
27     }
28     nr_free += n;
29     list_add(&free_list, &(base->page_link));
30 }
```

继续先看不认识的函数:

1. PageProperty 函数

写于: kern/mm/memlayout.h

```
1 | #define PageProperty(page)      test_bit(PG_property, &((page)->flags))
```

这个函数跟 **PageReserved** 函数是一个意思

不过这次是检查 (page)->flags 这个值的第 1 位是否为 0

如果为 0, 就可以顺利通过 assert(PageProperty(p)); 语句

如果为 1, 则报告错误并终止程序

说白了就是如果该页不是连续空闲空间的第一页, 那么就不报错, 顺利通过

2. list_next 函数

写于: libs/list.h

```
1 | /* *
2 |  * list_next - get the next entry
3 |  * @listelm:   the list head
4 |  */
5 | static inline list_entry_t *
6 | list_next(list_entry_t *listelm) {
7 |     return listelm->next;
8 | }
```

字面意思, 获得该列表的下一个元素

3. offsetof 函数

写于: libs/defs.h

```
1 | /* Return the offset of 'member' relative to the beginning of a struct
   | type */
2 | #define offsetof(type, member) \
3 |     ((size_t)(&((type *)0)->member))
```

又学到了, 先说 &((type *)0)->member 是啥:

ANSI C 标准允许值为 0 的常量被强制转换成任何一种类型的指针, 并且转换的结果是 NULL

因此 ((type *)0) 的结果就是一个类型为 **type *** 的 NULL 指针

如果利用这个 NULL 指针来访问 type 的成员当然是非法的

但 &((type *)0)->member 意图是想取 member 的地址

所以编译器不会生成访问 type 成员的代码, 而会优化为直接取地址

该函数的作用是获取结构体中 member 成员相对于该结构体首元素地址的偏移量

4. to_struct 函数

写于: libs/defs.h

```

1  /* *
2   * to_struct - get the struct from a ptr
3   * @ptr:      a struct pointer of member
4   * @type:     the type of the struct this is embedded in
5   * @member:   the name of the member within the struct
6   * */
7  #define to_struct(ptr, type, member) \
8      ((type *)((char *)(ptr) - offsetof(type, member)))

```

ptr 为一个真实地址，同时也是 **type *** 结构体中 member 成员变量的地址
 减去自身在 **type *** 中的偏移就可以得到这个结构体的起始地址
 那么这个函数就是用来获得这个成员变量所在结构体的起始地址

5. le2page 函数

写于: kern/mm/memlayout.h

```

1  // convert list entry to page
2  #define le2page(le, member) \
3      to_struct((le), struct Page, member)

```

经过 grep 命令，可以发现这个函数里的 le 变量基本都来自这里：**list_entry_t *le = list_next(&free_list);**

这个函数的作用就是依靠作为 Page 结构体中 member 成员变量的 le 变量，得到 le 成员变量所对应的结构体头变量

6. clear_bit 函数

写于: libs/atomic.h

```

1  /* *
2   * clear_bit - Atomically clears a bit in memory
3   * @nr:      the bit to clear
4   * @addr:    the address to start counting from
5   * */
6  static inline void
7  clear_bit(int nr, volatile void *addr) {
8      asm volatile ("btrl %1, %0" : "=m" (*(volatile long *)addr) : "Ir"
9      (nr));
10 }

```

btr 在执行 bt 命令的同时，把操作数的指定位置为 0
 所以这个函数是用来将 *addr 的第 nr 位的值设置为 0 的

7. ClearPageProperty 函数

写于: kern/mm/memlayout.h

```

1  #define ClearPageProperty(page)    clear_bit(PG_property, &((page)-
    >flags))

```

这个函数就是把 (page)->flags 的第 1 位设置为 0
 意思就是将该页标记为并非是连续空闲空间的第一页

8. __list_del 函数

写于: `libs/list.h`

```
1  /* *
2   * Delete a list entry by making the prev/next entries point to each
   other.
3   *
4   * This is only for internal list manipulation where we know
5   * the prev/next entries already!
6   */
7  static inline void
8  __list_del(list_entry_t *prev, list_entry_t *next) {
9      prev->next = next;
10     next->prev = prev;
11 }
```

双向链表删除元素的日常操作，就是把中间的元素断开，将上下两个元素相连
但是这里没有清除掉中间的元素，算是一个心里的疑惑

9. list_del 函数

写于: `libs/list.h`

```
1  /* *
2   * list_del - deletes entry from list
3   * @listelm:    the element to delete from the list
4   *
5   * Note: list_empty() on @listelm does not return true after this, the
   entry is
6   * in an undefined state.
7   */
8  static inline void
9  list_del(list_entry_t *listelm) {
10     __list_del(listelm->prev, listelm->next);
11 }
```

将 `listelm->prev` 与 `listelm->next` 相连，使 `listelm` 节点断连

Answer:

```
1  static void
2  default_free_pages(struct Page *base, size_t n) {
3      assert(n > 0); // 防止出现 n == 0 的情况发生
4      struct Page *p = base; // 设置 p 变量，一会要遍历，先从头开始
5      /* 开始遍历所有的物理空闲页，从 base 开始，经过 n 页后结束 */
6      for (; p != base + n; p++) {
7          /* 如果该页既是连续空闲空间的第一页，又被系统保留了，就不能通过 */
8          assert(!PageReserved(p) && !PageProperty(p));
9          p->flags = 0; // 将记录标志位的变量置零
10         set_page_ref(p, 0); // 将记录映射此物理页的虚拟页的变量的数值清零
11     }
12     base->property = n; // 将第一页结构体中记录连续空闲页数量的变量更新为 n
13     SetPageProperty(base); // 将第一页的 flags 变量中标记上该页是第一页
14     list_entry_t *le = list_next(&free_list); // 获取 &free_list 变量的下一个
   列表内容
15     /* *
16     * 循环遍历双向链表，初始地址为 &free_list 的下一个地址
```

```

17     * 直到遍历到 &free_list 地址，其实就是遍历整个链表
18     * 但是不进入链表头 &free_list 进行操作
19     * */
20     while (le != &free_list) {
21         p = le2page(le, page_link); // 获取当前链表块所对应的 Page 结构的结构体
头地址
22         /* *
23             下文没有用到 le 的，那么这里的意思就是更新 le 为下一个地址
24             为的是在进入下一个 while 语句后，寻找其所对应的 Page 结构体头
25             * */
26         le = list_next(le);
27         /* 被释放的空闲块后面紧接着空闲块 p */
28         if (base + base->property == p) {
29             /* *
30                 * 更新 base 页表头的 property 变量
31                 * 将接下来要被释放的页表头的 property 变量加到 base 页表头的 property
变量上
32                 * 差不多就是合并后面的空闲块 p 的意思
33                 * */
34             base->property += p->property;
35             /* *
36                 * 清除 p 结构体 flags 变量上的 Property 标志位
37                 * 意思就是让 p 不再成为某一页的页表头
38                 * */
39             ClearPageProperty(p);
40             list_del(&(p->page_link)); // 在链表中删除掉 p->page_link 的地址，
即释放掉 p
41         }
42         /* 被释放的空闲块前面紧接着空闲块 p */
43         else if (p + p->property == base) {
44             p->property += base->property; // 合并前面的空闲块 p
45             ClearPageProperty(base); // 让 base 不再成为某一页的页表头
46             base = p; // 更新 base 的地址为 p 的地址
47             list_del(&(p->page_link)); // 释放掉自己
48         }
49     }
50     nr_free += n; // 更新总空闲块数目
51     le = list_next(&free_list); // 获取 &free_list 的下一个地址
52     /* 与上面的 while 循环一样的遍历机制，不过这里新加了 break 的条件 */
53     while (le != &free_list) {
54         p = le2page(le, page_link); // 获取当前链表块所对应的 Page 结构的结构体
头地址
55         /* *
56             * 合并完成之后，由于需要维持链表有序，需要找到这个新的空闲块在链表中的插入位置
57             * 这里采用了插入排序的思想，遍历链表，找到第一个地址比需要插入的空闲块结尾地址大
的空闲块，插入到它的前面即可
58             * 由于前面的过程保证了不会有空闲块越界和粘连的情况发生，所以判断条件改为 base
<= p 也是可以的
59             * */
60             if (base + base->property <= p) {
61                 assert(base + base->property != p);
62                 break;
63             }
64             le = list_next(le);
65         }
66         list_add_before(le, &(base->page_link));
67     }

```

default_alloc_pages 函数

先看初始函数：

```
1 static struct Page *
2 default_alloc_pages(size_t n) {
3     assert(n > 0);
4     if (n > nr_free) {
5         return NULL;
6     }
7     struct Page *page = NULL;
8     list_entry_t *le = &free_list;
9     while ((le = list_next(le)) != &free_list) {
10         struct Page *p = le2page(le, page_link);
11         if (p->property >= n) {
12             page = p;
13             break;
14         }
15     }
16     if (page != NULL) {
17         list_del(&(page->page_link));
18         if (page->property > n) {
19             struct Page *p = page + n;
20             p->property = page->property - n;
21             list_add(&free_list, &(p->page_link));
22         }
23         nr_free -= n;
24         ClearPageProperty(page);
25     }
26     return page;
27 }
```

Answer:

```
1 static struct Page *
2 default_alloc_pages(size_t n) {
3     assert(n > 0); // 防止出现 n == 0 的情况发生
4     /* 若是申请的空闲块数目大于原有的空闲块数目，就返回 NULL */
5     if (n > nr_free) {
6         return NULL;
7     }
8     struct Page *page = NULL; // 声明一个 Page 结构体变量
9     list_entry_t *le = &free_list; // 声明一个 list_entry_t 结构体变量
10    // TODO: optimize (next-fit)
11    /* *
12     根据 First-Fit 算法，从地址小到大遍历每一个空闲块，一旦找到了足够大的空闲块就马上分配
13     也就是分配第一个足够大的空闲块给程序使用。找不到，说明没有这么大的连续空闲块，则返回空指针
14     * */
15    while ((le = list_next(le)) != &free_list) {
16        struct Page *p = le2page(le, page_link);
17        if (p->property >= n) {
18            page = p;
19            break;
20        }
21    }
```

```

21     }
22     /* *
23     如果找到了这样的空闲块，还要检查空闲块的大小是否比需要分配的大小更大
24     如果是，则需要进行分裂。分裂的操作便是将空闲块的后面的空闲块重新形成空闲块
25     这个空闲块的第一个页面是 page + n，大小是 page->property - n
26     设置好 property 和 flags 中的 property 位后，在链表中加入到当前空闲块的后面
    (因为地址更大)
27     */
28     if (page != NULL) {
29         if (page->property > n) {
30             struct Page *p = page + n;
31             p->property = page->property - n;
32             SetPageProperty(p);
33             list_add_after(&(page->page_link), &(p->page_link));
34         }
35         list_del(&(page->page_link)); // 删除 page 节点
36         nr_free -= n; // 更新物理空闲块数目
37         ClearPageProperty(page); // page 页不再是第一页
38     }
39     return page;
40 }

```

至此就完成了 First-Fit 算法

练习2：实现寻找虚拟地址对应的页表项

先看初始函数：

```

1 //get_pte - get pte and return the kernel virtual address of this pte for
  la
2 //          - if the PT contains this pte didn't exist, alloc a page for PT
3 // parameter:
4 // pgdir: the kernel virtual base address of PDT
5 // la: the linear address need to map
6 // create: a logical value to decide if alloc a page for PT
7 // return vaule: the kernel virtual address of this pte
8 pte_t *
9 get_pte(pde_t *pgdir, uintptr_t la, bool create) {
10     /* LAB2 EXERCISE 2: YOUR CODE
11     *
12     * If you need to visit a physical address, please use KADDR()
13     * please read pmm.h for useful macros
14     *
15     * Maybe you want help comment, BELOW comments can help you finish the
    code
16     *
17     * Some Useful MACROS and DEFINES, you can use them in below
    implementation.
18     * MACROS or Functions:
19     * PDX(la) = the index of page directory entry of VIRTUAL ADDRESS la.
20     * KADDR(pa) : takes a physical address and returns the corresponding
    kernel virtual address.
21     * set_page_ref(page,1) : means the page be referenced by one time
22     * page2pa(page): get the physical address of memory which this
    (struct Page *) page manages
23     * struct Page * alloc_page() : allocation a page

```

```

24      *   memset(void *s, char c, size_t n) : sets the first n bytes of the
memory area pointed by s
25      *
to the specified value c.
26      *   DEFINES:
27      *   PTE_P           0x001           // page table/directory
entry flags bit : Present
28      *   PTE_W           0x002           // page table/directory
entry flags bit : Writeable
29      *   PTE_U           0x004           // page table/directory
entry flags bit : User can access
30      */
31  #if 0
32      pde_t *pdep = NULL;    // (1) find page directory entry
33      if (0) {              // (2) check if entry is not present
34                          // (3) check if creating is needed, then alloc
page for page table
35                          // CAUTION: this page is used for page table, not
for common data page
36                          // (4) set page reference
37      uintptr_t pa = 0;    // (5) get linear address of page
38                          // (6) clear page content using memset
39                          // (7) set page directory entry's permission
40      }
41      return NULL;          // (8) return page table entry
42  #endif
43  }

```

先看不了解的定义和函数：

PDXSHIFT 宏

写于：kern/mm/mmu.h

```

1  #define PDXSHIFT      22           // offset of PDX in a linear
address

```

PDX 宏

写于：kern/mm/mmu.h

```

1  // page directory index
2  #define PDX(la) (((uintptr_t)(la)) >> PDXSHIFT) & 0x3FF

```

page_init 函数

写于：kern/mm/pmm.c

这里写了下面函数中 pages 的定义

```

1  /* pmm_init - initialize the physical memory management */
2  static void
3  page_init(void) {
4      struct e820map *memmap = (struct e820map *) (0x8000 + KERNBASE);
5      uint64_t maxpa = 0;
6

```



```

7     cprintf("e820map:\n");
8     int i;
9     for (i = 0; i < memmap->nr_map; i++) {
10        uint64_t begin = memmap->map[i].addr, end = begin + memmap-
>map[i].size;
11        cprintf("  memory: %08llx, [%08llx, %08llx], type = %d.\n",
12                memmap->map[i].size, begin, end - 1, memmap->map[i].type);
13        if (memmap->map[i].type == E820_ARM) {
14            if (maxpa < end && begin < KMEMSIZE) {
15                maxpa = end;
16            }
17        }
18    }
19    if (maxpa > KMEMSIZE) {
20        maxpa = KMEMSIZE;
21    }
22
23    extern char end[];
24
25    npage = maxpa / PGSIZE;
26    pages = (struct Page *)ROUNDUP((void *)end, PGSIZE);
27
28    for (i = 0; i < npage; i++) {
29        SetPageReserved(pages + i);
30    }
31
32    uintptr_t freemem = PADDR((uintptr_t)pages + sizeof(struct Page) *
npage);
33
34    for (i = 0; i < memmap->nr_map; i++) {
35        uint64_t begin = memmap->map[i].addr, end = begin + memmap-
>map[i].size;
36        if (memmap->map[i].type == E820_ARM) {
37            if (begin < freemem) {
38                begin = freemem;
39            }
40            if (end > KMEMSIZE) {
41                end = KMEMSIZE;
42            }
43            if (begin < end) {
44                begin = ROUNDUP(begin, PGSIZE);
45                end = ROUNDDOWN(end, PGSIZE);
46                if (begin < end) {
47                    init_memmap(pa2page(begin), (end - begin) / PGSIZE);
48                }
49            }
50        }
51    }
52 }

```

page2ppn 函数

写于: kern/mm/pmm.h

```

1 | static inline ppn_t
2 | page2ppn(struct Page *page) {
3 |     return page - pages;
4 | }

```

这个函数的意思是：**获取物理页对应的物理页号**

page2pa 函数

写于：kern/mm/pmm.h

```

1 | static inline uintptr_t
2 | page2pa(struct Page *page) {
3 |     return page2ppn(page) << PGSHIFT;
4 | }

```

PGSHIFT 的定义在 kern/mm/mmu.h 中，如下：

```

1 | #define PGSHIFT      12                // log2(PGSIZE)

```

该函数的意思是：**获取物理页对应的物理地址**

该地址向左移 12 位后就是 PPN，PPN 既是物理地址的高 20 位也是 PDE、PTE 的高 20 位

__panic 函数

写于：kern/debug/panic.c

```

1 | static bool is_panic = 0;
2 |
3 | /* *
4 |  * __panic - __panic is called on unresolvable fatal errors. it prints
5 |  * "panic: 'message'", and then enters the kernel monitor.
6 |  * */
7 | void
8 | __panic(const char *file, int line, const char *fmt, ...) {
9 |     if (is_panic) {
10 |         goto panic_dead;
11 |     }
12 |     is_panic = 1;
13 |
14 |     // print the 'message'
15 |     va_list ap;
16 |     va_start(ap, fmt);
17 |     cprintf("kernel panic at %s:%d:\n    ", file, line);
18 |     vcprintf(fmt, ap);
19 |     cprintf("\n");
20 |
21 |     cprintf("stack traceback:\n");
22 |     print_stackframe();
23 |
24 |     va_end(ap);
25 |
26 | panic_dead:
27 |     intr_disable();

```

```

28 |     while (1) {
29 |         kmonitor(NULL);
30 |     }
31 | }

```

该函数的作用在英文的已经有了标注：在不可解决的致命错误上调用。它输出 `panic: 'message'`，然后进入内核监视器。

panic 宏

写于：kern/debug/assert.h

```

1 | #define panic(...) \
2 |     __panic(__FILE__, __LINE__, __VA_ARGS__)

```

读代码就知道这个宏定义是对 `__panic` 函数的封装

PTXSHIFT 宏

写于：kern/mm/mmu.h

```

1 | #define PTXSHIFT      12                // offset of PTX in a linear
   | address

```

PPN 宏

写于：kern/mm/mmu.h

```

1 | // page number field of address
2 | #define PPN(la) (((uintptr_t)(la)) >> PTXSHIFT)

```

KADDR 宏

写于：kern/mm/pmm.h

```

1 | /* *
2 |  * KADDR - takes a physical address and returns the corresponding kernel
   | virtual
3 |  * address. It panics if you pass an invalid physical address.
4 |  * */
5 | #define KADDR(pa) ({ \
6 |     uintptr_t __m_pa = (pa); \
7 |     size_t __m_ppn = PPN(__m_pa); \
8 |     if (__m_ppn >= npage) { \
9 |         panic("KADDR called with invalid pa %08lx", __m_pa); \
10 |     } \
11 |     (void *) (__m_pa + KERNBASE); \
12 | })

```

看英文的意思，该宏定义的作用就是通过物理地址找到对应的逻辑(虚拟)地址

PDE_ADDR 宏

写于：kern/mm/mmu.h

```

1 #define PTE_ADDR(pte) ((uintptr_t)(pte) & ~0xFFF)
2 #define PDE_ADDR(pde) PTE_ADDR(pde)

```

Answer:

```

1 pte_t *
2 get_pte(pde_t *pgdir, uintptr_t la, bool create) {
3     /* *
4      * pdep -- Page Directory Entry Pointer
5      * pgdir: the kernel virtual base address of PDT
6      * PDX(la) = the index of page directory entry of VIRTUAL ADDRESS la.
7      * 这里的 pgdir 可以是页目录表的基址
8      * PDX 的作用是找到线性地址 la 中负责页目录表的高 10 位的对应下标，正好是左移 22
    位得到
9      * 然后利用 pgdir，根据下标找到 PDE，里面的高 20 位存着 Page Table 的基址右移
    12 位的值
10     * 将其赋给指针变量 pdep
11     * */
12     pde_t *pdep = &pgdir[PDX(la)];
13     /* *
14     * 检查查找到的页目录项是否存在，如果存在直接放回找到的页表项即可//
15     * 即检查是否设置了 Present 位，也就是 PDE_P 位
16     * 不过实际上并没有 PDE_P 这个宏，使用注释里告诉我们用等价的 PTE_P 来检查
17     * 注释里告诉我们这个位是 PDE 和 PTE 通用的
18     * */
19     if (!(*pdep & PTE_P)) {
20         struct Page *page;
21         /* *
22         * 如果该页目录项是不存在的，或者物理空间不足 (page)，那么就不创建新的页表
    (!create)，直接返回 NULL
23         * alloc_pages 函数就是之前练习一里写的 default_alloc_pages 函数
24         * 通过 alloc_pages() 分配的页地址 并不是真正的页分配的地址，实际上只是 Page
    这个结构体所在的地址而已
25         * 故而需要通过使用 page2pa() 将 Page 这个结构体的地址转换成真正物理页地址的
    线性地址
26         * 然后需要注意的是无论是 * 或是 memset 都是对虚拟地址进行操作的
27         * 所以需要将真正的物理页地址再转换成内核逻辑(虚拟)地址
28         * */
29         if (!create || (page = alloc_page()) == NULL) {
30             return NULL;
31         }
32         set_page_ref(page, 1); // 更新该物理页的引用计数
33         uintptr_t pa = page2pa(page); // 获取物理页对应的物理地址的基址，即左移
    12 位
34         /* *
35         * KADDR 的作用是将该地址加上内核基地址得到该项正确的逻辑地址，也就是得到了页表
    (PTT)真正的逻辑地址
36         * 然后对新创建的页表进行初始化，赋予可被用户访问、可写权限及存在标志
37         * */
38         memset(KADDR(pa), 0, PGSIZE);
39         *pdep = pa | PTE_U | PTE_W | PTE_P; // 设置 PTE 的权限
40     }
41     /* *
42     * PDE_ADDR(*pdep) 清空了 *pdep 的后 12 位，这个值是页表 (PTT) 的偏移地址
43     * KADDR 的作用是将该地址加上内核基地址得到该项正确的逻辑地址，也就是得到了页表
    (PTT)真正的逻辑地址

```

```

44      * PTX 的作用是找到线性地址 la 中负责页表的中 10 位的对应下标，正好是左移 12 位得
    到
45      * 得到的就是 PTE 的地址，这是个指针，所以最后返回的值就是 ptep
46      * */
47      return &((pte_t *)KADDR(PDE_ADDR(*pdep)))[PTX(la)];
48  }

```

练习3：释放某虚地址所在的页并取消对应二级页表项的映射

源码：

```

1  //page_remove_pte - free an Page struct which is related linear address la
2  //                  - and clean(invalidate) pte which is related linear
    address la
3  //note: PT is changed, so the TLB need to be invalidate
4  static inline void
5  page_remove_pte(pde_t *pgdir, uintptr_t la, pte_t *ptep) {
6      /* LAB2 EXERCISE 3: YOUR CODE
7       *
8       * Please check if ptep is valid, and tlb must be manually updated if
    mapping is updated
9       *
10      * Maybe you want help comment, BELOW comments can help you finish the
    code
11      *
12      * Some Useful MACROS and DEFINES, you can use them in below
    implementation.
13      * MACROS or Functions:
14      *   struct Page *page pte2page(*ptep): get the according page from the
    value of a ptep
15      *   free_page : free a page
16      *   page_ref_dec(page) : decrease page->ref. NOTICE: ff page->ref == 0
    , then this page should be free.
17      *   tlb_invalidate(pde_t *pgdir, uintptr_t la) : Invalidate a TLB
    entry, but only if the page tables being
18      *                  edited are the ones currently in use by the
    processor.
19      * DEFINES:
20      *   PTE_P           0x001           // page table/directory
    entry flags bit : Present
21      */
22  #if 0
23      if (0) { // (1) check if this page table entry is
    present
24          struct Page *page = NULL; // (2) find corresponding page to pte
25                                     // (3) decrease page reference
26                                     // (4) and free this page when page
    reference reaches 0
27                                     // (5) clear second page table entry
28                                     // (6) flush tlb
29      }
30  #endif
31  }

```

page_ref_dec 函数

写于: kern/mm/pmm.h

```
1 static inline int
2 page_ref_dec(struct Page *page) {
3     page->ref -= 1;
4     return page->ref;
5 }
```

将该页的引用计数减 1

free_page 函数

写于: kern/mm/pmm.h

```
1 #define free_page(page) free_pages(page, 1)
```

这个就是练习一里咱们自己写的 **default_free_pages** 函数

rcr3 函数

写于: libs/x86.h

```
1 static inline uintptr_t
2 rcr3(void) {
3     uintptr_t cr3;
4     asm volatile ("mov %%cr3, %0" : "=r" (cr3) :: "memory");
5     return cr3;
6 }
```

将 cr3 寄存器的值赋给 cr3 变量, 返回 cr3 变量的值

PADDR 宏

写于: kern/mm/pmm.h

```
1 /* *
2  * PADDR - takes a kernel virtual address (an address that points above
3  * KERNBASE),
4  * where the machine's maximum 256MB of physical memory is mapped and
5  * returns the
6  * corresponding physical address. It panics if you pass it a non-kernel
7  * virtual address.
8  * */
9 #define PADDR(kva) ({
10     uintptr_t __m_kva = (uintptr_t)(kva);
11     if (__m_kva < KERNBASE) {
12         panic("PADDR called with invalid kva %08lx", __m_kva);
13     }
14     __m_kva - KERNBASE;
15 })
```

PADDR -- 接受内核虚拟地址(位于 KERNBASE 之上的地址), 映射机器的最大 256MB 物理内存, 并返回相应的物理地址

如果你给它一个非内核虚拟地址, 它会进入 panic 函数; 也就是说, 这个宏跟 **KADDR** 宏是对着来的

正常情况下, 这个函数的作用就是**将当前的虚拟地址减去 KERNBASE 来得到对应的物理地址**

invlpg 函数

写于: `libs/x86.h`

```
1 static inline void
2 invlpg(void *addr) {
3     asm volatile ("invlpg (%0)" :: "r" (addr) : "memory");
4 }
```

补补知识。。俺就看啥都是第一次见

实际的地址转换过程中, 处理器首先依据线性地址查找 TLB, 假设未发现该线性地址到物理地址的映射关系 (TLB miss)

将依据页表中的映射关系填充 TLB (TLB fill), 然后再进行地址转换

该指令使与按地址指向的内存关联的页面的翻译后备链表缓冲区 (TLB) 失效。

tlb_invalidate 函数

写于: `kern/mm/pmm.c`

```
1 // invalidate a TLB entry, but only if the page tables being
2 // edited are the ones currently in use by the processor.
3 void
4 tlb_invalidate(pde_t *pgdir, uintptr_t la) {
5     if (rcr3() == PADDR(pgdir)) {
6         invlpg((void *)la);
7     }
8 }
```

判断 cr3 寄存器里的值是否等于 PDT 的物理地址, 如果等于那么就取消 la 对应物理页之间的关联, 相当于刷新 TLB

Answer:

```
1 static inline void
2 page_remove_pte(pde_t *pgdir, uintptr_t la, pte_t *ptep) {
3     if ((*ptep & PTE_P)) {
4         struct Page *page = pte2page(*ptep);
5         if (page_ref_dec(page) == 0) { // 若引用计数减一后为 0, 则释放该物理页
6             free_page(page);
7         }
8         *ptep = 0; // 清空 PTE
9         tlb_invalidate(pgdir, la); // 刷新 TLB
10    }
11 }
```

扩展练习Challenge: buddy system (伙伴系统) 分配算法

以后再来

扩展练习Challenge：任意大小的内存单元slub分配算法

以后再来

参考链接：

https://blog.csdn.net/weixin_34327223/article/details/85450897

https://blog.csdn.net/abraham_li/article/details/49474721

<https://blog.csdn.net/rosetta/article/details/90746936>

<http://ilinuxkernel.com/?p=1062>