

OS lab2

成员信息

解子萱 2312585 信息安全、法学双学位班

崔颖欣 2311136 信息安全、法学双学位班

范鼎辉 2312326 信息安全

练习1：理解first-fit 连续物理内存分配算法（思考题）

first-fit 连续物理内存分配算法作为物理内存分配一个很基础的方法，需要同学们理解它的实现过程。请大家仔细阅读实验手册的教程并结合 `kern/mm/default_pmm.c` 中的相关代码，认真分析 `default_init`, `default_init_memmap`, `default_alloc_pages`, `default_free_pages`等相关函数，并描述程序在进行物理内存分配的过程以及各个函数的作用。请在实验报告中简要说明你的设计实现过程。请回答如下问题：

- 你的first fit算法是否有进一步的改进空间？

准备工作

在开始解析first-fit 连续物理内存分配算法之前，有必要对一些关键头文件进行整理，理解内存的基本实现结构和基本原理设计。

空闲页链表的实现：通过list.h 定义环形双向链表结构

基本双向链表实现如下：

结构体/函数	代码实现	功能解析
<code>list_entry</code> 节点结构体	<pre>struct list_entry { struct list_entry *prev, *next; };</pre>	<ul style="list-style-type: none">每个节点的 <code>prev</code> 指针指向前一个节点；<code>next</code> 指针指向下一个节点<code>typedef</code> 定义 <code>list_entry_t</code> 为 <code>struct list_entry</code> 的别名
<code>list_init</code> 链表初始化	<pre>static inline void list_init(list_entry_t *elm) {</pre>	<ul style="list-style-type: none">传入链表节点，将该节点的 <code>prev</code> 和 <code>next</code> 都指向自己，作为孤立的表头

	<pre>elm->prev = elm->next = elm; }</pre>	
__list_add 节点插入的底层函数	<pre>static inline void __list_add(list_entry_t *elm, list_entry_t *prev, list_entry_t *next) { prev->next = next->prev = elm; elm->next = next; elm->prev = prev; }</pre>	<ul style="list-style-type: none"> 将新节点 <code>elm</code> 插入到已知的前一个节点 <code>prev</code> 和下一个节点 <code>next</code> 之间，并调整它们的 <code>next</code> 和 <code>prev</code> 指针。
__list_del 节点删除的底层函数	<pre>static inline void __list_del(list_entry_t *prev, list_entry_t *next) { prev->next = next; next->prev = prev; }</pre>	<ul style="list-style-type: none"> 删除节点 <code>prev</code> 和 <code>next</code> 之间的节点。“前的后为后，后的前为前”，跳过中间的节点，实现删除操作。
list_empty 检查链表是否为空	<pre>static inline bool list_empty(list_entry_t *list) {return list->next == list; }</pre>	<ul style="list-style-type: none"> 如果链表为空，则 <code>list->next</code> 会指向 <code>list</code> 本身，即初始化时的状态（孤立表头）
list_next 获取下一个节点	<pre>static inline list_entry_t *list_next(list_entry_t *listelm) {return listelm->next; }</pre>	<ul style="list-style-type: none"> 通过 <code>listelm->next</code> 获取下一个节点的指针。
list_prev 获取上一个节点	<pre>static inline list_entry_t *list_prev(list_entry_t *listelm) {return listelm->prev; }</pre>	<ul style="list-style-type: none"> 通过 <code>listelm->prev</code> 获取前一个节点的指针。

在此基础上，一些包装后的上层函数包括：

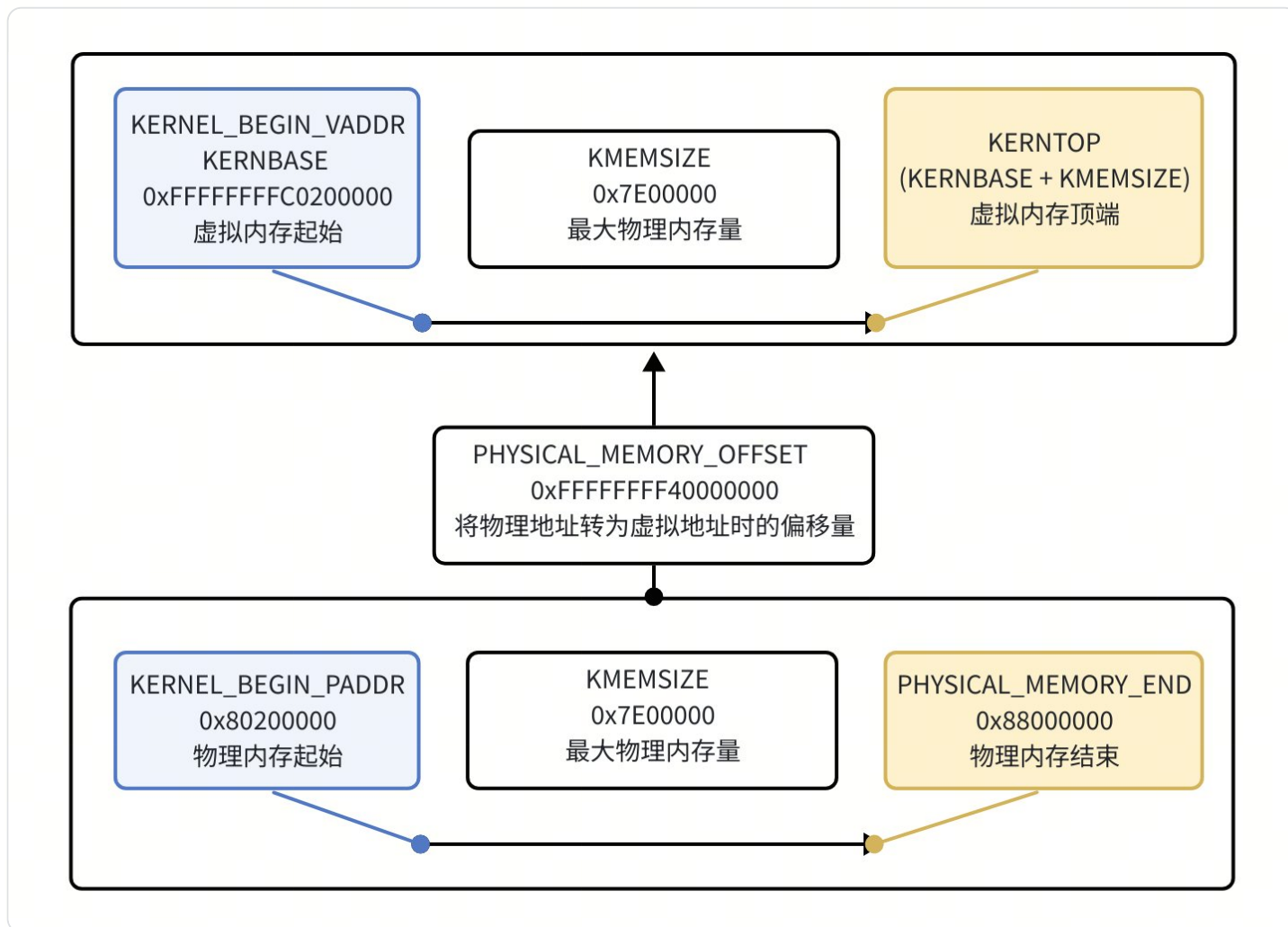
- `list_add` 节点插入封装（默认向后）：调用 `list_add_after`，将新节点 `elm` 添加到旧节点 `listelm` 之后
- `list_add_before` 节点向前插入的封装：调用 `__list_add`，将新节点 `elm` 在 `listelm` 节点之前插入。
- `list_add_after` 节点向后插入的封装：调用 `__list_add`，将新节点 `elm` 在 `listelm` 节点之后插入。

- `list_del` 节点删除封装：调用 `__list_del`，将 `listelm` 节点从链表中移除
- `list_del_init` 节点删除+初始化封装：先调用 `list_del` 删除节点，再调用 `list_init` 初始化节点，使其成为一个孤立的节点。

内存结构实现：通过memlayout.h 定义内存布局结构

顶部地址映射宏定义

关键示意图总结如下：



其中`0x80000000`到`0x80200000`被OpenSBI占用

栈页大小宏定义

主要涉及下面两句代码：

```
#define KSTACKPAGE      2
#define KSTACKSIZE      (KSTACKPAGE * PGSIZE)
```

- 意味着内核栈由 2 页组成，栈大小为2*页大小

Page页结构定义

根据代码定义，一个页结构拥有以下属性：

- `ref` 引用计数，`ref == 0` 表示当前没有引用
- `flags` 标志位，用来表示页面的各种状态，两个标志位（bit）：

名称	含义	位编号
PG_reserved	页被内核保留，用户分配器不能使用；若 PG_reserved = 1，代表这个页不可分配；若 PG_reserved = 0，代表这个页可以被分配或管理。	第 0 位（最低位）
PG_property	页是空闲块的头页（管理连续空闲页块）；如果一个空闲块包含连续的多页，只有头页会设置 PG_property；	第 1 位

- `property` 连续空闲块的页数：当一个连续空闲区以某页为头时，该页的 `property` 等于该连续空闲块的页数，而非头页则把 `property` 设为 0。
- `page_link` 一个 `list_entry_t` 链表节点，用于将此 `Page` 节点插入到空闲链表 `free_list` 中

特殊宏——`le2page` 宏：已知 `Page` 中的 `page_link` 字段的地址，把它转回到对应的 `Page` 指针

Free_area空闲区结构定义

- `free_area_t` 作为struct结构，其实是封装保存一个空闲页链表的头节点（`free_list`）和该链表中页的总数（`nr_free`）。

default_pmm.c 代码分析

default_init 初始化工作

我们首先针对全局内存空区进行了一系列初始化准备工作：

代码块

```
1 //内存空闲区定义，
2 static free_area_t free_area;
3 #define free_list (free_area.free_list)
4 #define nr_free (free_area.nr_free)
5 //初始化函数，初始化空页链表的头节点，设置总节点（空页表）数量为0：
6 static void
7 default_init(void) {
8     list_init(&free_list);
9     nr_free = 0;
10 }
11 //获取空页表数量的函数
```

```

12  static size_t
13  default_nr_free_pages(void) {
14      return nr_free;
15  }

```

default_init_memmap 空闲页初始化工作

这一部分我们将连续的物理页（从 base 开始共 n 个）初始化为可分配的空闲页：

代码块

```

1  static void
2  default_init_memmap(struct Page *base, size_t n) {
3      assert(n > 0);
4      struct Page *p = base;
5      //遍历每一页，确保：1.都是保留页 2.清零页属性 3.清零页引用计数
6      for (; p != base + n; p++) {
7          assert(PageReserved(p));
8          p->flags = p->property = 0;
9          set_page_ref(p, 0);
10     }
11     //设置头页“base”的属性，并把这 n 页加入空闲总量统计
12     base->property = n;
13     SetPageProperty(base);
14     nr_free += n;
15     //如果当前空闲链表为空，就直接把这块插入链表头后面
16     if (list_empty(&free_list)) {
17         list_add(&free_list, &(base->page_link));
18     }
19     //否则，遍历整个空闲链表，寻找合适的插入位置：
20     else {
21         list_entry_t* le = &free_list;
22         while ((le = list_next(le)) != &free_list) {
23             struct Page* page = le2page(le, page_link);
24             //如果 base < page: 说明新块地址更靠前，按地址顺序插入到 page 前面；
25             if (base < page) {
26                 list_add_before(le, &(base->page_link));
27                 break;
28             }
29             //如果已经到达链表尾，表示新块地址比所有已有块都大，把它加在链表尾部。
30             else if (list_next(le) == &free_list) {
31                 list_add(le, &(base->page_link));
32             }
33         }
34     }
35 }

```

default_alloc_pages 内存页 First Fit 分配算法

其主体思想是：在空闲页链表 free_list 中找到第一块连续的空闲页，并从中划出 n 页供使用。

代码块

```
1  static struct Page *
2  default_alloc_pages(size_t n) {
3      //大小检查
4      assert(n > 0);
5      if (n > nr_free) {
6          return NULL;
7      }
8      //遍历空闲链表，寻找合适的空闲页page
9      struct Page *page = NULL;
10     list_entry_t *le = &free_list;
11     while ((le = list_next(le)) != &free_list) {
12         struct Page *p = le2page(le, page_link);
13         // 关键点: p->property 表示该空闲块的页数，找到第一个 property >= n 的块，即
        First Fit;
14         // 若找到，则保存到 page 并 break。
15         if (p->property >= n) {
16             page = p;
17             break;
18         }
19     }
20     // 如果找到了合适的page，进行以下操作：
21     if (page != NULL) {
22         //先取出当前块的前驱结点 prev，再将该空闲块从空闲链表中删除（因为要分配它）；
23         list_entry_t* prev = list_prev(&(page->page_link));
24         list_del(&(page->page_link));
25         //处理分配后的剩余部分
26         if (page->property > n) {
27             struct Page *p = page + n;           //新的空闲块头页 = page + n;
28             p->property = page->property - n; //它的 property = 原property - n;
29             SetPageProperty(p);                //设置 PG_property 位;
30             list_add(prev, &(p->page_link)); //把它插在原空闲块的前驱后面，保持地址
        顺序
31         }
32         // 减少系统空闲页计数；并清除 PG_property，标记这块页已不再是空闲块；
33         nr_free -= n;
34         ClearPageProperty(page);
35     }
36     return page;
37 }
```

default_free_pages 物理页释放算法

其主要思想是：释放从 base 开始的 n 个连续页Page，并将它们加入系统的空闲页链表 free_list 中

代码块

```
1  static void
2  default_free_pages(struct Page *base, size_t n) {
3      assert(n > 0);
4      struct Page *p = base;
5      //对每一页检查它：1.不是保留页 2.不属于已标记的空闲块；
6      //随后清空标志位与引用计数，防止被错误认为仍在使用的
7      for (; p != base + n; p++) {
8          assert(!PageReserved(p) && !PageProperty(p));
9          p->flags = 0;
10         set_page_ref(p, 0);
11     }
12     //设定空闲块头信息
13     base->property = n;
14     SetPageProperty(base);
15     nr_free += n;
16     //将获得的空闲页插入 free_list，并保证按物理地址升序排列，逻辑与初始化时一致，故不
    再重复
17     if (list_empty(&free_list)) {
18         list_add(&free_list, &(base->page_link));
19     }
20     else {
21         list_entry_t* le = &free_list;
22         while ((le = list_next(le)) != &free_list) {
23             struct Page* page = le2page(le, page_link);
24             if (base < page) {
25                 list_add_before(le, &(base->page_link));
26                 break;
27             } else if (list_next(le) == &free_list) {
28                 list_add(le, &(base->page_link));
29             }
30         }
31     }
32     // 空闲页向前合并的步骤
33     list_entry_t* le = list_prev(&(base->page_link));
34     if (le != &free_list) {
35         p = le2page(le, page_link);
36         // 看链表中 base 前一个空闲块 p 是否紧邻
37         //条件 p + p->property == base表示物理上连续，若相邻则合并：
38         if (p + p->property == base) {
39             p->property += base->property;
40             ClearPageProperty(base);
41             list_del(&(base->page_link));
42             //关键点：更新新的块头是更前面的页
```

```

43         base = p;
44     }
45 }
46 // 空闲页向后合并的步骤，与上一致，不再重复
47 le = list_next(&(base->page_link));
48 if (le != &free_list) {
49     p = le2page(le, page_link);
50     if (base + base->property == p) {
51         base->property += p->property;
52         ClearPageProperty(p);
53         list_del(&(p->page_link));
54     }
55 }
56 }

```

check 内存设置自检测函数

basic_check代码及检测总结如下：

代码块

```

1  static void
2  basic_check(void) {
3      struct Page *p0, *p1, *p2;
4      p0 = p1 = p2 = NULL;
5      assert((p0 = alloc_page()) != NULL);
6      assert((p1 = alloc_page()) != NULL);
7      assert((p2 = alloc_page()) != NULL);
8
9      assert(p0 != p1 && p0 != p2 && p1 != p2);
10     assert(page_ref(p0) == 0 && page_ref(p1) == 0 && page_ref(p2) == 0);
11
12     assert(page2pa(p0) < npage * PGSIZE);
13     assert(page2pa(p1) < npage * PGSIZE);
14     assert(page2pa(p2) < npage * PGSIZE);
15
16     list_entry_t free_list_store = free_list;
17     list_init(&free_list);
18     assert(list_empty(&free_list));
19
20     unsigned int nr_free_store = nr_free;
21     nr_free = 0;
22
23     assert(alloc_page() == NULL);
24
25     free_page(p0);
26     free_page(p1);

```



```

27     free_page(p2);
28     assert(nr_free == 3);
29
30     assert((p0 = alloc_page()) != NULL);
31     assert((p1 = alloc_page()) != NULL);
32     assert((p2 = alloc_page()) != NULL);
33
34     assert(alloc_page() == NULL);
35
36     free_page(p0);
37     assert(!list_empty(&free_list));
38
39     struct Page *p;
40     assert((p = alloc_page()) == p0);
41     assert(alloc_page() == NULL);
42
43     assert(nr_free == 0);
44     free_list = free_list_store;
45     nr_free = nr_free_store;
46
47     free_page(p);
48     free_page(p1);
49     free_page(p2);
50 }

```

阶段	动作	检查点	意图
1	分配三页	页不重复、ref=0	验证 alloc 基本正确性
2	检查物理地址	地址在合法范围内	验证 page→pa 映射
3	清空空闲链表	alloc 应失败	测试极端条件
4	释放页后再分配	成功分配3页	检查 free/alloc 互通
5	释放单页再分配	返回相同页	验证链表与计数逻辑
6	恢复链表状态	无副作用	保持全局一致性
7	释放残留页	内存清理	防止污染全局环境

default_check代码及检测总结如下：

代码块

```

1     static void
2     default_check(void) {
3         int count = 0, total = 0;
4         list_entry_t *le = &free_list;
5         while ((le = list_next(le)) != &free_list) {
6             struct Page *p = le2page(le, page_link);

```

```
7         assert(PageProperty(p));
8         count ++, total += p->property;
9     }
10    assert(total == nr_free_pages());
11
12    basic_check();
13
14    struct Page *p0 = alloc_pages(5), *p1, *p2;
15    assert(p0 != NULL);
16    assert(!PageProperty(p0));
17
18    list_entry_t free_list_store = free_list;
19    list_init(&free_list);
20    assert(list_empty(&free_list));
21    assert(alloc_page() == NULL);
22
23    unsigned int nr_free_store = nr_free;
24    nr_free = 0;
25
26    free_pages(p0 + 2, 3);
27    assert(alloc_pages(4) == NULL);
28    assert(PageProperty(p0 + 2) && p0[2].property == 3);
29    assert((p1 = alloc_pages(3)) != NULL);
30    assert(alloc_page() == NULL);
31    assert(p0 + 2 == p1);
32
33    p2 = p0 + 1;
34    free_page(p0);
35    free_pages(p1, 3);
36    assert(PageProperty(p0) && p0->property == 1);
37    assert(PageProperty(p1) && p1->property == 3);
38
39    assert((p0 = alloc_page()) == p2 - 1);
40    free_page(p0);
41    assert((p0 = alloc_pages(2)) == p2 + 1);
42
43    free_pages(p0, 2);
44    free_page(p2);
45
46    assert((p0 = alloc_pages(5)) != NULL);
47    assert(alloc_page() == NULL);
48
49    assert(nr_free == 0);
50    nr_free = nr_free_store;
51
52    free_list = free_list_store;
53    free_pages(p0, 5);
```

```

54
55     le = &free_list;
56     while ((le = list_next(le)) != &free_list) {
57         struct Page *p = le2page(le, page_link);
58         count --, total -= p->property;
59     }
60     assert(count == 0);
61     assert(total == 0);
62 }

```

阶段	检查项目	验证点
1	链表遍历	所有空闲块有 PageProperty
2	全局计数	nr_free_pages() 与链表统计一致
3	分配单块	能成功分配连续页块
4	分配失败	无空闲页时 alloc_page() 返回 NULL
5	部分释放	释放中间页段后能正确形成空闲块
6	块拆分	分配部分页后剩余块保留正确 property
7	块合并	释放相邻块后能自动合并
8	链表有序	分配时按地址顺序维护
9	计数一致性	所有操作后 count/total 不变

default_pmm_manager物理内存分配器初始化

这个结构体本质上是一个函数指针表，描述该物理内存管理器应当实现的所有接口函数，相当于对上述所有工作做整合封装。

代码块

```

1  const struct pmm_manager default_pmm_manager = {
2      .name = "default_pmm_manager",      // 初始化内存管理器的名字
3      .init = default_init,               // 初始化函数（如初始化链表、锁等）
4      .init_memmap = default_init_memmap, // 初始化内存页结构（把物理页标记为空闲）
5      .alloc_pages = default_alloc_pages, // 分配指定数量的物理页
6      .free_pages = default_free_pages,   // 释放物理页
7      .nr_free_pages = default_nr_free_pages, // 查询当前空闲页数
8      .check = default_check,             // 检查函数，用于调试或验证正确性
9  };

```

程序物理内存分配过程

下面总结 `default_pmm.c` 在整个物理内存分配中的作用：

- `default_init` 初始化整个全局空页链表，并设置总节点数量为0；
- `default_init_memmap` 将连续的n个物理页初始化为可分配的空闲页，并插入空闲链表；
- `default_alloc_pages` 实现FirstFit内存分配算法，在空闲页链表中找到第一块连续的空闲页，并从中划出 n 页供使用；
- `default_free_pages` 释放 n 个连续的内存页，并将它们加入空闲页链表中，随后完成空页上下文合并。

总的来说，系统启动时，`default_init` 建立物理页管理的基础结构；

`default_init_memmap` 形成可供分配的内存池。当内核或程序请求分配内存时，`default_alloc_pages` 按照First Fit算法将所需页划出并标记为已用；当内存被释放时，`default_free_pages` 负责将这些页重新插入空闲链表，并自动与相邻空闲块合并。整个流程实现了物理页从初始化、分配到回收的完整管理机制。

算法改进空间

1. 对于free_list的遍历性能优化

当前 `default_init_memmap`、`default_alloc_pages` 和 `default_free_pages` 都是从头到尾对空闲页链表进行遍历：

代码块

```
1  list_entry_t* le = &free_list;
2  while ((le = list_next(le)) != &free_list) {
3      ...
4  }
```

很明显这是线性 $O(n)$ 查找插入点，如果空闲页块较多，每次释放或初始化都会线性扫描，性能会急剧下降！因此可以引入一个简单的按地址有序插入二叉搜索树，来把释放或插入复杂度降为 $O(\log n)$ 。

2. 空闲页合并逻辑优化

在 `default_free_pages()` 中源代码分别作了前向和后向的合并，但是这两个操作逻辑是一样的，因此我们其实可以把他们合并封装在一起，成为一个merge函数，使得逻辑语义更清晰。

代码块

```
1  //输入: Page *left, Page *right
2  //处理:
3      if (left + left->property == right) {
4          left->property += right->property;
```

```
5         ClearPageProperty(right);
6         list_del(&(right->page_link));
7     }
8     //使用：分别执行两次该函数：merge(prev_page, base) 和 merge(base, next_page)
```

练习2：实现 Best-Fit 连续物理内存分配算法（需要编程）

在完成练习一后，参考kern/mm/default_pmm.c对First Fit算法的实现，编程实现Best Fit页面分配算法，算法的时空复杂度不做要求，能通过测试即可。请在实验报告中简要说明你的设计实现过程，阐述代码是如何对物理内存进行分配和释放，并回答如下问题：

- 你的 Best-Fit 算法是否有进一步的改进空间？

First Fit /Best Fit 算法对比：

最先匹配法First-fit

- 按分区的先后次序，从头查找，找到符合要求的第一个分区就分配。
- 分配和释放的时间性能较好，较大的空闲分区可以被保留在内存高端。但随着低端分区不断划分而产生较多小分区，每次分配时查找时间开销会增大。

最佳匹配法Best-fit

- 将分区按小大顺序组织,找到的第一个适应分区是大小与要求相差最小的空闲分区。
- 整体来看会形成较多外碎片。但较大的空闲分区可以被保留。

Best Fit 页面分配算法实现：

逻辑关键

Best Fit 关键在于如何找到**大小最合适**的那一块内存

设计思路

对于页面分配来说，我们选择引入一个额外参数min_size，初始化为比空闲页总数还大的值；之后在对空闲页链表进行遍历，记录实时大小“最合适”的页并更新min_size，直到最后遍历结束。而对于页面释放来说，与之前First Fit算法相同，故不再赘述。

代码实现

Best Fit仅仅**只需要修改空闲页选取的逻辑**，其余代码皆和First Fit算法一样，主要修改的地方在下方标黄代码处：

代码块

```
1    //在空闲页链表 free_list 中找到一块大小最合适的连续空闲页，并从中划出n页供使用。
```

```

2  static struct Page *
3  Best_Fit_alloc_pages(size_t n) {
4      //依然大小检查
5      assert(n > 0);
6      if (n > nr_free) {
7          return NULL;
8      }
9      //遍历空闲链表, 寻找大小最合适的空闲页page
10     size_t min_size = nr_free + 1; //记录最合适的大小, 这里初始化为能表示的最大值
11     struct Page *page = NULL;
12     list_entry_t *le = &free_list;
13     while ((le = list_next(le)) != &free_list) {
14         struct Page *p = le2page(le, page_link);
15         // 如果有大小合适的空闲页
16         if (p->property >= n) {
17             //如果比当前最合适的页还要合适, 就选取他
18             if (p->property - n < min_size) {
19                 min_size = p->property - n;
20                 page = p;
21             }
22         }
23     }
24     // 如果找到了合适的page, 进行以下操作:
25     if (page != NULL) {
26         //先取出当前块的前驱结点 prev, 再将该空闲块从空闲链表中删除 (因为要分配它);
27         list_entry_t* prev = list_prev(&(page->page_link));
28         list_del(&(page->page_link));
29         //处理分配后的剩余部分
30         if (page->property > n) {
31             struct Page *p = page + n;           //新的空闲块头页 = page + n;
32             p->property = page->property - n; //它的 property = 原property - n;
33             SetPageProperty(p);                 //设置 PG_property 位;
34             list_add(prev, &(p->page_link)); //把它插在原空闲块的前驱后面, 保持地址
35             //顺序
36             }
37             // 减少系统空闲页计数; 并清除 PG_property, 标记这块页已不再是空闲块;
38             nr_free -= n;
39             ClearPageProperty(page);
40         }
41         return page;
42     }

```

测试结果

将测试目标在pmm.c中改为best_fit算法之后, 调用make grade, 获取执行结果

lab2

可以看到，输出测试结果为Best Fit——OK，证明我们的算法编写正确。

算法改进空间

与First Fit算法代码一样，它对于空闲页链表的遍历仍为线性 $O(n)$ 查找，如果空闲页块较多，每次释放或初始化都会线性扫描，性能会急剧下降，故我们还是可以引入一个简单的按地址有序插入二叉搜索树，来把释放或插入复杂度降为 $O(\log n)$ 。

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

Buddy System算法把系统中的可用存储空间划分为存储块(Block)来进行管理, 每个存储块的大小必须是2的n次幂(2^n), 即1, 2, 4, 8, 16, 32, 64, 128...

- 参考[伙伴分配器的一个极简实现](#)，在ucore中实现buddy system分配算法，要求有比较充分的测试用例说明实现的正确性，需要有设计文档。

buddy system设计文档

1 算法简介

Buddy System 是一种将可用内存划分为大小为 2 的幂次的块来进行管理的分配算法。我们在分配时总是选择能够覆盖请求的最小幂次块，若块比请求大则递归地分裂为“伙伴”块，直到得到合适大小；在释放时与其伙伴检测并尽量合并成更大的块，以减少外部碎片。这套机制用一棵完全二叉树来表达每个节点对应内存区域的最大可用块大小，根节点覆盖整个区域，叶节点对应最小单位页。

本实现参考了wuwenbin的buddy system极简实现，但根据ucore的pmm框架进行了适配和重构，并在一些边界处理和接口兼容方面进行了优化，使其能够无缝集成到ucore操作系统的物理内存管理体系中。在实现顺序上，我们先定义数据结构和树索引宏，然后实现初始化，再完成分配/释放逻辑，最后提供查询与校验。

2 核心数据结构

我们用一个结构体来描述 Buddy System 的核心状态：总页数（2 的幂次）、树数组（每个节点最大可用块大小）、物理页基址。**它借鉴了wuwenbin实现中的二叉树思想，但适配了ucore的Page**

结构。

代码块

```
1 struct buddy_system {
2     unsigned size;           // 管理的总页面数，必须是2的幂
3     unsigned *longest;      // 二叉树数组，记录每个节点对应的最大可用块大小
4     struct Page *base;      // 内存区域的起始页面
5 };
```

这里的 `size` 是逻辑管理规模（向上取整到 2 的幂），`longest` 是长度为 `2*size-1` 的完全二叉树数组，`base` 指向 ucore 的 `Page` 数组的起始位置。树的叶子对应最小页，父节点记录左右子树的最大可用块，以便快速定位分配位置。

3 关键宏定义

为了简化代码实现，我们借鉴了wuwenbin实现中的一些关键宏定义，用简单宏来表达树的结构关系与常用判断，便于实现算法时保持代码直观清晰。

代码块

```
1 #define LEFT_LEAF(index) ((index) * 2 + 1)
2 #define RIGHT_LEAF(index) ((index) * 2 + 2)
3 #define PARENT(index) (((index) + 1) / 2 - 1)
4 #define IS_POWER_OF_2(x) (!(x) & ((x) - 1))
5 #define MAX(a, b) ((a) > (b) ? (a) : (b))
```

`LEFT_LEAF`/`RIGHT_LEAF`/`PARENT` 表示在完全二叉树上的索引关系，不需要额外指针；`IS_POWER_OF_2` 用位运算判断幂次；`MAX` 是父节点汇聚左右子状态时的取值规则。

4 pmm_manager结构体

我们在 .c 文件末尾将 buddy System 的各个操作函数以 `pmm_manager` 结构体的方式对外暴露。正如实验课中宫老师所讲，操作系统设计者用了较为“tricky”的方式，让 C 也能实现类似 C++ 类实例化的效果：每种新的物理内存分配算法在 .h 中都声明一个 `pmm_manager`，在 .c 末尾都实例化一个对象，像是在实例化一个“类”的对象一样。这使得替换不同算法只需切换这个实例，接口统一、扩展方便。

代码块

```
1 const struct pmm_manager buddy_system_pmm_manager = {
2     .name = "buddy_system_pmm_manager",
3     .init = buddy_init,
4     .init_memmap = buddy_init_memmap,
5     .alloc_pages = buddy_alloc_pages,
```



```
6     .free_pages = buddy_free_pages,
7     .nr_free_pages = buddy_nr_free_pages,
8     .check = buddy_check,
9 };
```

5 核心函数

接下来，我们展示每个函数关键算法对应的代码片段。需要说明的是，由于本部分代码量较大，为节省篇幅，不再给出完整代码实现。完整代码及详细注释可参见项目的 [GitHub 仓库](#)。

5.1 buddy_init：初始化状态占位

`buddy_init` 函数负责初始化Buddy System的数据结构。这个函数的设计思路是先进行简单的初始化，实际的详细初始化工作放在 `buddy_init_memmap` 中进行：

代码块

```
1  static void
2  buddy_init(void) {
3      buddy.size = 0;
4      buddy.longest = NULL;
5      buddy.base = NULL;
6  }
```

5.2 buddy_init_memmap：构建树、标记可用叶子、处理非幂次规模

这个函数是整个Buddy System的初始化核心，它根据给定的内存区域初始化Buddy System。我们需要先计算合适的树大小，并初始化所有节点的值。关键算法部分包括：

代码块

```
1  // size向上取整到2的幂，并限制最大size防止数组越界
2  unsigned size = 1; while (size < n) size <<= 1;
3  if (size > 2048) size = 2048;
4
5  static unsigned tree[4096];           // 2*2048-1 = 4095
6  buddy.longest = tree; buddy.size = size; buddy.base = base;
7  memset(tree, 0, sizeof(tree));
8
9  unsigned node_size = size * 2;
10 for (int i = 0; i < 2 * size - 1; ++i) {
11     if (IS_POWER_OF_2(i + 1)) node_size /= 2;
12     buddy.longest[i] = node_size;
13 }
14
15 // 屏蔽超出实际页面数的叶子，并自底向上更新父节点
```

```

16  if (n < size) {
17      for (unsigned i = n; i < size; i++) {
18          unsigned index = i + size - 1;
19          buddy.longest[index] = 0;
20      }
21      for (int i = size - 2; i >= 0; i--) {
22          unsigned left = buddy.longest[LEFT_LEAF(i)];
23          unsigned right = buddy.longest[RIGHT_LEAF(i)];
24          buddy.longest[i] = MAX(left, right);
25      }
26  }

```

这部分代码体现了Buddy System的核心设计思想：将内存划分为2的幂次大小的块，并使用完全二叉树来管理这些内存块。

5.3 buddy_alloc_pages：向下查找合适节点、标记分配、向上汇总

分配指定数量的页面是这个算法的核心功能。我们从根节点开始，递归查找合适大小的空闲块，找到后标记为已分配，并更新父节点的信息。关键算法代码：

代码块

```

1  // 请求大小向上取整为2的幂
2  unsigned size = 1; while (size < n) size <<= 1;
3
4  // 从根向下选择能容纳size的子树
5  unsigned index = 0, node_size;
6  if (buddy.longest[index] < size) return NULL;
7  for (node_size = buddy.size; node_size != size; node_size /= 2) {
8      if (buddy.longest[LEFT_LEAF(index)] >= size) index = LEFT_LEAF(index);
9      else index = RIGHT_LEAF(index);
10 }
11
12 // 标记叶子已分配并计算页偏移
13 buddy.longest[index] = 0;
14 unsigned offset = (index + 1) * node_size - buddy.size;
15
16 // 向上更新父节点的最大可用块大小
17 while (index) {
18     index = PARENT(index);
19     buddy.longest[index] = MAX(buddy.longest[LEFT_LEAF(index)],
20                               buddy.longest[RIGHT_LEAF(index)]);
21 }

```

这个算法展示了Buddy System的高效性，通过二叉树结构快速定位到合适的内存块。

5.4 buddy_free_pages: 定位节点、标记空闲、伙伴合并

释放指定的页面时，我们需要找到对应的节点，标记为空闲，然后检查伙伴节点是否也空闲，如果是则合并。关键算法：

代码块

```
1 // 根据页偏移定位叶子索引
2 unsigned offset = base - buddy.base;
3 unsigned node_size = 1, index = offset + buddy.size - 1;
4
5 // 向上查找“值为0的节点”，找到真实已分配块的叶子
6 for (; buddy.longest[index]; index = PARENT(index)) {
7     node_size *= 2;
8     if (index == 0) return; // 防御
9 }
10
11 // 标记空闲并尝试与伙伴合并
12 buddy.longest[index] = node_size;
13 while (index) {
14     index = PARENT(index); node_size *= 2;
15     unsigned L = buddy.longest[LEFT_LEAF(index)];
16     unsigned R = buddy.longest[RIGHT_LEAF(index)];
17     buddy.longest[index] = (L + R == node_size) ? node_size : MAX(L, R);
18 }
```

5.5 buddy_nr_free_pages: 查询当前最大可分配块

这个函数返回当前可用的页面数，实现非常简单：

代码块

```
1 static size_t
2 buddy_nr_free_pages(void) {
3     if (buddy.longest == NULL) return 0;
4     return buddy.longest[0];
5 }
```

6 测试用例

测试部分包含两个函数：`basic_check` 和 `buddy_check`。`basic_check` 我们直接复用了 `default_pmm.c` 的实现，而 `buddy_check` 则是我们重点设计的专项测试。

6.1 basic_check

`basic_check` 函数直接复用了 `default_pmm.c` 中的实现，这是为了确保Buddy System能够通过ucore框架的基础测试。这个测试函数主要验证了基本的页面分配和释放功能，包括：

- 分配单个页面并验证分配成功
- 验证分配的页面地址在合理范围内
- 验证页面引用计数正确
- 测试页面释放功能
- 验证内存泄漏检查

这个基础测试确保了Buddy System与ucore框架的兼容性，为我们后续的专项测试打下了坚实的基础。

6.2 buddy_check函数

`buddy_check` 函数是我们专门为Buddy System设计的综合性测试函数，包含了200多行代码，全面测试了Buddy System的各种功能和边界情况。测试内容主要包括：

- **测试1：基本的2的幂次分配和释放** - 验证Buddy System能够正确分配和释放1、2、4、8页等2的幂次大小的内存块。
- **测试2：非2的幂次分配（向上取整）** - 测试请求3、5、6、7页时，系统能够正确向上取整到4、8、8、8页进行分配。
- **测试3：伙伴合并机制测试** - 验证当连续分配多个小内存块并释放时，系统能够正确合并伙伴块形成更大的空闲块。
- **测试4：连续分配和释放压力测试** - 通过连续分配和释放16个1页块，测试系统在压力情况下的稳定性和性能。
- **测试5：大块分配测试** - 测试分配16页大块的能力，验证系统对大内存请求的处理。
- **测试6：边界情况测试** - 包括测试分配1024页和512页等边界情况，验证系统对极端内存请求的处理。
- **测试7：碎片整理效果测试** - 通过有策略地分配和释放内存块，制造碎片化场景，然后测试系统是否仍能分配大块内存。
- **测试8：功能正确性验证** - 验证分配的页面地址是否在合理范围内，分配的页面地址是否唯一，确保没有地址冲突。
- **测试9：最终内存泄漏检查** - 比较测试前后的可用页面数，确保没有内存泄漏发生。

这些测试用例覆盖全面，输出带有明确的“成功/失败/警告”提示，基本上把Buddy System核心行为和边界都跑了一遍，且在QEMU下全部通过，确保了算法的正确性和健壮性。最终测试结果展示如下：

```

=== 开始buddy system专项测试 ===
初始可用页面数：2048

--- 测试1：基本的2的幂次分配和释放 ---
✓ 分配1页成功
✓ 分配2页成功
✓ 分配4页成功
✓ 分配8页成功
分配后可用页面数：1024
✓ 释放1页成功
✓ 释放2页成功
✓ 释放4页成功
✓ 释放8页成功

--- 测试2：非2的幂次分配（向上取整） ---
✓ 请求3页分配成功（实际分配4页）
✓ 请求5页分配成功（实际分配8页）
✓ 请求6页分配成功（实际分配8页）
✓ 请求7页分配成功（实际分配8页）
✓ 释放3页成功
✓ 释放5页成功
✓ 释放6页成功
✓ 释放7页成功

--- 测试3：伙伴合并机制测试 ---
✓ 分配第1个1页块成功
✓ 分配第2个1页块成功
✓ 分配第3个1页块成功
✓ 分配第4个1页块成功
合并前可用页面数：1024
✓ 释放第1个1页块
✓ 释放第2个1页块
✓ 释放第3个1页块
✓ 释放第4个1页块
合并后可用页面数：2048

--- 测试4：连续分配和释放压力测试 ---
✓ 成功分配16个1页块
✓ 释放所有分配的块

--- 测试5：大块分配测试 ---
✓ 成功分配第1个16页大块
✓ 成功分配第2个16页大块
✓ 成功分配第3个16页大块
✓ 成功分配第4个16页大块
✓ 释放第1个16页大块
✓ 释放第2个16页大块
✓ 释放第3个16页大块
✓ 释放第4个16页大块

--- 测试6：边界情况测试 ---
✓ 分配1024页成功
✓ 释放1024页成功
✓ 分配512页成功
✓ 释放512页成功

--- 测试7：碎片整理效果测试 ---
✓ 分配第1个2页块
✓ 分配第2个2页块
✓ 分配第3个2页块
✓ 分配第4个2页块
✓ 分配第5个2页块
✓ 分配第6个2页块
✓ 分配第7个2页块
✓ 分配第8个2页块
✓ 释放第2个2页块（制造碎片）
✓ 释放第4个2页块（制造碎片）
✓ 释放第6个2页块（制造碎片）
✓ 释放第8个2页块（制造碎片）
✓ 碎片化后仍能分配8页大块

--- 最终检查 ---
最终可用页面数：2048
✓ 内存泄漏检查通过

--- 测试8：功能正确性验证 ---
✓ 分配的页面地址在合理范围内

```

答辩后补充：

先前的测试函数只显示分配是否成功，但没有显示具体分配的内存地址和范围，无法直观验证buddy system的正确性，因此我们需要优化测试函数。首先添加辅助函数 `show_page_info`：

代码块

```

1 // 辅助函数：显示被分配页面的详细地址信息
2 static void show_page_info(struct Page *page, size_t n, const char *desc) {
3     // 计算页面在pages数组中的索引
4     size_t page_index = page - buddy.base; // 页面索引 = 页面地址 - 基址
5     // 使用系统提供的page2pa函数计算物理地址
6     uintptr_t phys_addr = page2pa(page); // 调用系统函数获取物理地址
7
8     cprintf(" %s: 页面索引[%d-%d], 物理地址[0x%08x-0x%08x], 大小%dKB\n",
9             desc, page_index, page_index + n - 1, // 页面索引范围
10            phys_addr, phys_addr + n * PGSIZE - 1, // 物理地址范围
11            n * PGSIZE / 1024); // 大小 (KB)
12 }

```

该函数能显示被分配页面的详细地址信息。接着，我们在测试函数中调用辅助函数，例如：

代码块

```

1 p1 = alloc_pages(1); // 分配1页
2 assert(p1 != NULL); // 确保分配成功
3 cprintf("✓ 分配1页成功\n"); // 输出成功信息
4 show_page_info(p1, 1, "1页块"); // 调用辅助函数

```


这样就可以输出被分配页面的具体地址信息，以验证buddy system的正确性。我们节选测试1结果如下：

```
--- 测试1: 基本的2的幂次分配和释放 ---
✓ 分配1页成功
  1页块: 页面索引[0-0], 物理地址[0x8034b000-0x8034bfff], 大小4KB
✓ 分配2页成功
  2页块: 页面索引[2-3], 物理地址[0x8034d000-0x8034efff], 大小8KB
✓ 分配4页成功
  4页块: 页面索引[4-7], 物理地址[0x8034f000-0x80352fff], 大小16KB
✓ 分配8页成功
  8页块: 页面索引[8-15], 物理地址[0x80353000-0x8035afff], 大小32KB
```

可以看出，输出结果的格式均为 页面索引[begin-end]，物理地址[begin-end]，大小xKB，并且测试1的结果是完全正确的。最初我们分配了[0]处的一个页，当我们再次需要两个页时，跳过了[1]处的页，而去分配[2-3]的两个页。后面的分配结果同理，不再赘述。我们再对其他测试结果进行节选展示如下：

```
--- 测试2: 非2的幂次分配（向上取整） ---
✓ 请求3页分配成功（实际分配4页）
  请求3页实际分配4页: 页面索引[0-3], 物理地址[0x8034b000-0x8034efff], 大小16KB
✓ 请求5页分配成功（实际分配8页）
  请求5页实际分配8页: 页面索引[8-15], 物理地址[0x80353000-0x8035afff], 大小32KB
✓ 请求6页分配成功（实际分配8页）
  请求6页实际分配8页: 页面索引[16-23], 物理地址[0x8035b000-0x80362fff], 大小32KB
✓ 请求7页分配成功（实际分配8页）
  请求7页实际分配8页: 页面索引[24-31], 物理地址[0x80363000-0x8036afff], 大小32KB
```

```
--- 测试3: 伙伴合并机制测试 ---
分配4个1页块，观察buddy system的分配策略：
✓ 分配第1个1页块成功
  1页块: 页面索引[0-0], 物理地址[0x8034b000-0x8034bfff], 大小4KB
✓ 分配第2个1页块成功
  1页块: 页面索引[1-1], 物理地址[0x8034c000-0x8034cfff], 大小4KB
✓ 分配第3个1页块成功
  1页块: 页面索引[2-2], 物理地址[0x8034d000-0x8034dfff], 大小4KB
✓ 分配第4个1页块成功
  1页块: 页面索引[3-3], 物理地址[0x8034e000-0x8034efff], 大小4KB
```

```
--- 测试5: 大块分配测试 ---
尝试分配多个16页大块，观察buddy system对大内存的处理：
✓ 成功分配第1个16页大块
  16页大块: 页面索引[0-15], 物理地址[0x8034b000-0x8035afff], 大小64KB
✓ 成功分配第2个16页大块
  16页大块: 页面索引[16-31], 物理地址[0x8035b000-0x8036afff], 大小64KB
✓ 成功分配第3个16页大块
  16页大块: 页面索引[32-47], 物理地址[0x8036b000-0x8037afff], 大小64KB
✓ 成功分配第4个16页大块
  16页大块: 页面索引[48-63], 物理地址[0x8037b000-0x8038afff], 大小64KB
```

至此，测试函数就完美辣！

7 总结

总体上，这份设计以“先有数据结构与树宏，才能有初始化；先有初始化，才能保障分配/释放逻辑正确”为主线，算法核心参考了 wuwenbin 的极简实现，但我们在 ucore 的接口规范、边界安全（数组越界与非幂次规模）、测试充分性等方面做了落地适配。接口层通过 `pmm_manager` 的“类实例化”式结构统一暴露，既达到了替换简单、扩展友好的目标，也与课程框架保持一致。

扩展练习Challenge：任意大小的内存单元slub分配算法（需要编程）

slub算法，实现两层架构的高效内存单元分配，第一层是基于页大小的内存分配，第二层是在第一层基础上实现基于任意大小的内存分配。可简化实现，能够体现其主体思想即可。

- 参考[linux的slub分配算法/](#)，在ucore中实现slub分配算法。要求有比较充分的测试用例说明实现的正确性，需要有设计文档。

SLUB设计文档

SLAB vs SLUB：

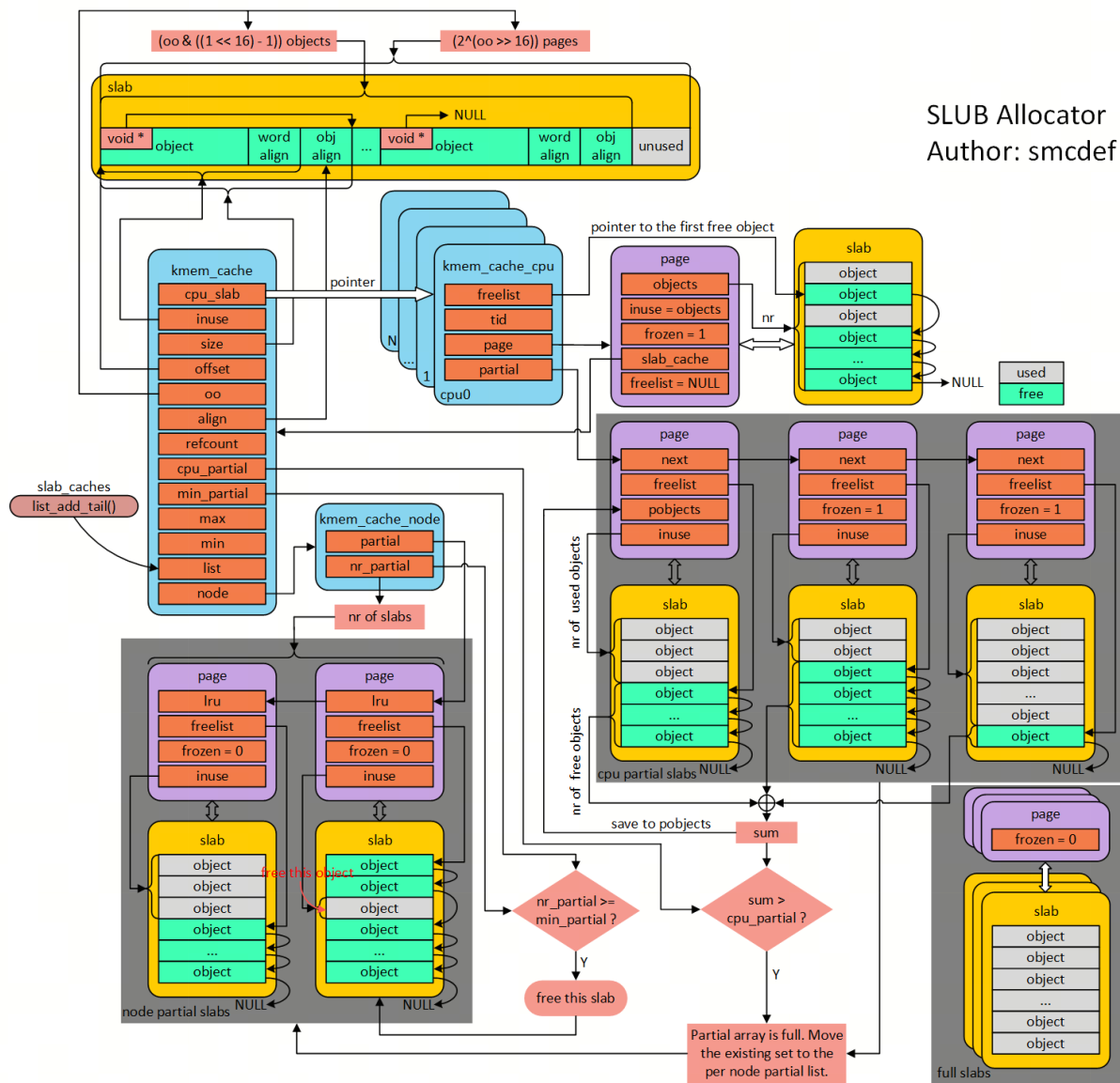
早期 SLAB 通过每个 `kmem_cache` 维护 full/partial/empty 队列与 bufctl/bitmap 记录空闲对象，多核下容易产生锁竞争。SLUB（unqueued）**取消了 SLAB 风格的中心化对象队列与 bufctl，改用对象内嵌 freelist**，并通过每 CPU 当前活动 slab 的快路径降低共享结构争用；慢路径才访问全局/NUMA 级别的 partial 列表，结构更简、扩展性更好。

本实验实现的是简化版 SLUB：保留 `kmem_cache`（多种 size）与 `slab`（按页切对象），使用对象内嵌 freelist 实现 O(1) 取还；未实现 Linux 真 SLUB 的 per-CPU 快路径与完整慢路径，但核心思路（两层架构 + 轻量元数据 + 大对象回退页级）已具备。

文中为避免歧义：“slab”指从 buddy system 拿来的页块（或页组）；“SLUB”指算法/分配器；代码里 `slub_` 前缀仅表示“本 SLUB 实现的函数或数据结构”。

1 算法简介

SLUB 是一种面向“小对象”的高效内存分配器。它并不直接管理整块物理内存，而是**站在页级分配器（如 Buddy System）肩膀上工作**：先从底层**按页**拿到内存，再把页**切成均匀小对象**，用简洁的空闲链来快速分配与回收。



我的实现是一个简化版 SLUB，重点体现两层架构与核心思路：

- **第一层：页级分配** → 全部委托给已实现的 Buddy System (`alloc_pages` / `free_pages`)
- **第二层：对象级分配** → 在每个 slab（若干页，简化为 1 或 2 页）里按固定对象大小切块，用 `freelist` 管空闲对象；小对象走 SLUB，大对象直接回退到页级。

在 uCore 中，物理内存管理模块通过统一的 `pmm_manager` 接口选择具体实现；我新增了 `slub_pmm_manager` 并在 `pmm.c` 中切换为默认管理器以运行测试。

2 实现思路

2.1 两层架构总览

chatgpt画的好看的图

- 1 应用/内核对象
- 2 | (按对象大小)
- 3 SLUB 层 (对象级)
- 4 |
- 5 | `kmem_cache(size=N)` | ← 多个缓存池 (32/64/.../2048B)


```

6   |   |— slabs_partial |   ← 还有空位的 slab
7   |   |— slabs_full   |   ← 满的 slab
8   |_____|
9       | (按页拿/还)
10  Buddy System 层 (页级)
11  |_____|
12  | alloc_pages/free... |
13  |_____|

```

- **小对象**：落到合适的 `kmem_cache`，从其 `slab` 的 `freelist` 取块；
- **大对象 ($\geq \text{PGSIZE}$)**：直接走 `alloc_pages`（页级回退），避免 SLUB 头部占用导致放不下的问题。

在我的实现里，`2048B` 为了更接近真实 SLUB 的“每 slab 容量更合理”，特地演示了**双页 slab** 的配置，其余 size class 多为单页 slab。

2.2 核心数据结构

文中为避免歧义：“slab”指从 buddy system 拿来的页块（或页组）；“SLUB”指算法/分配器；代码里 `slub_` 前缀仅表示“本 SLUB 实现的函数或数据结构”。

1. `kmem_cache`：每种对象大小一个缓存池

关键信息：`objsize`、`slabs_partial`、`slabs_full`、`slab_pages`（本实验中 1 或 2）。

我在 `slub_init_caches()` 里初始化了一组常见的 size class：`{32, 64, 128, 256, 512, 1024, 2048, 4096}`；并示范把 `2048B` 设置为双页 slab。

代码块

```

1  typedef struct kmem_cache {
2      const char *name;           // 方便调试的名字
3      unsigned int objsize;       // 对象大小
4      list_entry_t slabs_partial; // 还有空位的 slab 列表
5      list_entry_t slabs_full;    // 已满的 slab 列表
6      unsigned int slab_pages;    // 每个 slab 占用的页数（简化为1）
7  } kmem_cache_t;

```

2. `slub_slab_t`：一个 slab 对应一片连续页

保存这片页的起始 `Page *`、`inuse` 已用计数、`capacity` 总对象数、对象空闲链 `free_list` 等；

代码块

```

1  typedef struct slub_slab {

```

```

2     struct Page *page;           // 这个 slab 对应的起始页
3     list_entry_t link;           // 链接到 cache 的 slab 列表
4     list_entry_t free_list;      // slab 内部的空闲对象链表
5     unsigned int inuse;          // 当前已分配对象数
6     unsigned int objsize;        // 对象大小 (字节)
7     unsigned int capacity;       // 这个 slab 可以容纳的对象总数
8 } slub_slab_t;

```

简化点：我把 slab 元数据直接放在 slab 页的页头，真实的实现更精细一些，这样方便调试

3. obj_node：对象里的小链表头

直接用对象开头的若干字节当 `list_entry_t`，把空闲对象串起来，`0(1)` 取还。

代码块

```

1 // 把对象指针挂到 list_entry_t: 用对象地址作为 list_entry_t 存储，这里简单把对象首地址
  当成一个 list 节点
2 typedef struct obj_node {
3     list_entry_t link; // 链表节点在对象头 (简化: 对象就是裸内存, 前几个字节当节点用)
4 } obj_node_t;

```

2.2 核心函数

2.3.1 slub初始化

- `slub_init()`：初始化全部 `kmem_cache` (size class、链表、`slab_pages`)。
- `slub_init_memmap(base, n)`：直接把页级初始化转给 Buddy System，因为 SLUB 不负责底层。

代码块

```

1 // 初始化所有 size class 的 cache
2 static void slub_init_caches(void) {
3     for (unsigned i = 0; i < SLUB_SIZE_CLASS_COUNT; i++) {
4         caches[i].name = "slub";
5         caches[i].objsize = slub_sizes[i];
6         list_init(&caches[i].slabs_partial);
7         list_init(&caches[i].slabs_full);
8         caches[i].slab_pages = (slub_sizes[i] == 2048) ? 2 : 1; // 2048B 使用双
        页 slab
9     }
10 }

```

2.3.2 选 cache: `slub_select_cache(size)`

- 按“**最小可容纳**”原则选择 size class；超出上限则返回 `NULL`，触发页级回退；
- 我保留了 `4096` 这一档用于演示“正好一页”的情况（但由于有 slab 元数据，4096B 实际不走对象级而会回退页级，这里在分配时也做了兜底处理）。

代码块

```
1 // 根据 size 选择合适的 cache (向上取最近的 size class)
2 static kmem_cache_t *slub_select_cache(size_t size) {
3     for (unsigned i = 0; i < SLUB_SIZE_CLASS_COUNT; i++) {
4         if (size <= slub_sizes[i]) return &cache[i];
5     }
6     return NULL; // 超过 4096B 的对象不适用 SLUB (交给页级)
7 }
```

2.3.3 创建 slab: `slab_create(cache)`

- 向 Buddy System 申请 `slab_pages` 页；
- 页头写入 `slub_slab_t` 元数据；
- 计算 `capacity = (total_bytes - sizeof(slab_header)) / objsize`；
- 把每个对象用 `obj_node` 链进 `free_list`。

若 `capacity == 0`（比如对象过大），实际系统里要**回退页级或拒绝**；我的实现里在对象分配路径已有回退逻辑覆盖。

代码块

```
1 // 初始化一个 slab, 来自 buddy 系统分配的一页
2 static slub_slab_t *slab_create(kmem_cache_t *cache) {
3     struct Page *page = buddy_system_pmm_manager.alloc_pages(cache-
4 >slab_pages);
5     if (page == NULL) return NULL;
6     // 为 slab 元数据使用这页的起始地址的一小块空间 (不再另开页, 教学用简化)
7     void *page_va = (void *) (page2pa(page) + va_pa_offset);
8     memset(page_va, 0, PGSIZE * cache->slab_pages);
9
10    slub_slab_t *slab = (slub_slab_t *)page_va; // slab 元数据放在页头
11    slab->page = page;
12    list_init(&slab->link);
13    list_init(&slab->free_list);
14    slab->inuse = 0;
15    slab->objsize = cache->objsize;
16    unsigned total = PGSIZE * cache->slab_pages;
17    unsigned usable = total - sizeof(slub_slab_t);
18    slab->capacity = usable / cache->objsize;
```

```

19
20 // 切分对象，并把空闲对象都串起来（对象紧随 slab 元数据之后）
21 // 注意：这是教学简化，真实 SLUB 会把元数据放在 per-slab 外部或者在对象末尾等，不影响主体思想
22 unsigned char *obj_base = (unsigned char *)page_va + sizeof(slub_slab_t);
23
24 for (unsigned i = 0; i < slab->capacity; i++) {
25     obj_node_t *node = (obj_node_t *) (obj_base + i * cache->objsize);
26     list_add(&slab->free_list, &node->link);
27 }
28
29 return slab;
30 }

```

2.3.4 分配对象：slub_alloc(size)

- 先 `slub_select_cache` 找到 cache；若找不到或对象超大，**直接页级回退**（向 Buddy 要整页/多页，返回页的 VA）；
- 否则在 `slabs_partial` 里找第一个能分配出的 slab，走 `slab_alloc_obj` 取一个对象；
- 若 slab 用满了，把 slab 从 partial 移到 full；
- 如果 partial 里都塞满了，那就 `slab_create()` 一个新 slab，再分配。

代码块

```

1 // 提供一个对象分配接口（教学扩展），不改变 pmm_manager 的签名，测试中直接调用
2 static void *slub_alloc(size_t size) {
3     kmem_cache_t *cache = slub_select_cache(size);
4     if (cache == NULL || cache->objsize == PGSIZE) {
5         // 对象太大，直接分配整页（或多页）。这里返回页的虚拟地址（教学简化）。
6         struct Page *pg = buddy_system_pmm_manager.alloc_pages((size + PGSIZE
7 - 1)/PGSIZE);
8         if (!pg) return NULL;
9         return (void *) (page2pa(pg) + va_pa_offset);
10    }
11
12    // 尝试在 partial 列表找到空位对象
13    list_entry_t *le = &cache->slabs_partial;
14    while ((le = list_next(le)) != &cache->slabs_partial) {
15        // 教学简化：把 list 节点本身当做指向 slab 的节点，改为容器获取
16        // 我们在创建 slab 时使用 slab->link 作为链入节点，因此这里需要用 to_struct
17        slub_slab_t *slab = to_struct(le, slub_slab_t, link);
18        void *obj = slab_alloc_obj(slab);
19        if (obj != NULL) {
20            // 如果分配后满了，移到 full 列表
21            if (slab->inuse == slab->capacity) {

```

```

21         list_del(&slab->link);
22         list_add(&cache->slabs_full, &slab->link);
23     }
24     return obj;
25 }
26 }
27
28 // 没有可用对象就新建一个 slab
29 slub_slab_t *slab = slab_create(cache);
30 if (slab == NULL) return NULL;
31 // 新 slab 加入 partial 列表
32 list_add(&cache->slabs_partial, &slab->link);
33 void *obj = slab_alloc_obj(slab);
34 if (obj == NULL) return NULL;
35 if (slab->inuse == slab->capacity) {
36     list_del(&slab->link);
37     list_add(&cache->slabs_full, &slab->link);
38 }
39 return obj;
40 }

```

2.3.5 释放对象：slub_free(ptr, size)

- 通过对象地址反推所属 slab 的页起始（因为 slab header 在页头/双页头，很好定位）；
- 把对象挂回 free_list，inuse--；
- 若刚释放使得 slab 从满变非满，把 slab 从 full 移回 partial；
- （简化取舍）暂不做 slab 变空后的“整页回收”，这样可以减少抖动、便于复用（真实系统会有回收策略，这里留作改进方向）。

代码块

```

1 static void slub_free(void *obj, size_t size) {
2     kmem_cache_t *cache = slub_select_cache(size);
3     if (cache == NULL) {
4         // 当 size>4096 或自定义大对象，调用方需使用 slub_free_pages；这里无法安全定
        位页
5         return;
6     }
7     // 根据对象地址定位所属 slab：对象所在页的开头保存了 slab 元数据
8     uintptr_t obj_pa = (uintptr_t)obj - va_pa_offset;
9     uintptr_t slab_base_pa = ROUNDDOWN(obj_pa, PGSIZE * cache->slab_pages);
10    slub_slab_t *slab = (slub_slab_t *) (slab_base_pa + va_pa_offset);
11    slab_free_obj(slab, obj);
12    // 如果 slab 从 full 变成 partial，需要移动列表
13    if (slab->inuse + 1 == slab->capacity) {

```

```
14         // 原来是满的，刚释放后不满
15         list_del(&slab->link);
16         list_add(&cache->slabs_partial, &slab->link);
17     }
18     // 当 slab 变为空，可以考虑回收页：这里教学简化不做回收，保留以降低碎片重新分配成本
19 }
```

2.3.6 页级接口：slub_alloc_pages / slub_free_pages / slub_nr_free_pages

- 这三个接口**全部直通 Buddy System**，保持与 pmm_manager 其它实现兼容。SLUB 只接管“小对象”的那部分逻辑。

2.3.7 对接 pmm_manager & 启动路径

- 新管理器在 slub_pmm.c 末尾注册为：

```
const struct pmm_manager slub_pmm_manager = {... .check =
slub_check, ...};
```
- pmm.c 中选择并打印 memory management: slub_pmm_manager；调用 pmm_manager->init() / init_memmap() / check()；
- pmm.h 中 pmm_manager 接口定义见头文件。

3 测试

3.1 基础函数 basic_check()

用 Buddy 连续分配/释放 1 页，确认底层页级 OK；

3.2 自测函数 slub_check() 及测试结果

1. 各对象大小的 Cache 验证

这一段依次对 8 种对象大小（32B、64B、128B、256B、512B、1024B、2048B、4096B）进行了完整的分配与释放流程验证。

```

[SLUB] === slub_check() begin ===
[SLUB] basic_check() passed.
[SLUB] size=32 selected cache: objsize=32, slab_pages=1
[SLUB] new cache test for size=32 begin
[SLUB] first alloc: va=ffffffffc034dfd8 pa=8034dfd8 page_pa=8034d000 | partial=1 full=0
[SLUB] capacity per slab = 126 objects
[SLUB] alloc i=63: va=ffffffffc034d7f8 pa=8034d7f8 page_pa=8034d000 | partial=1 full=0
[SLUB] alloc i=125: va=ffffffffc034d038 pa=8034d038 page_pa=8034d000 | partial=0 full=1
[SLUB] after fill: partial=0 full=1
[SLUB] extra alloc: va=ffffffffc034efd8 pa=8034efd8 page_pa=8034e000 | first_page_pa=8034d000
[SLUB] after free first: partial=2 full=0
[SLUB] before count(partial) after free
[SLUB] partial count=2
[SLUB] size=32 test done.
[SLUB] size=64 selected cache: objsize=64, slab_pages=1
[SLUB] new cache test for size=64 begin
[SLUB] first alloc: va=ffffffffc034ffb8 pa=8034ffb8 page_pa=8034f000 | partial=1 full=0
[SLUB] capacity per slab = 63 objects
[SLUB] alloc i=31: va=ffffffffc034f7f8 pa=8034f7f8 page_pa=8034f000 | partial=1 full=0
[SLUB] alloc i=62: va=ffffffffc034f038 pa=8034f038 page_pa=8034f000 | partial=0 full=1
[SLUB] after fill: partial=0 full=1
[SLUB] extra alloc: va=ffffffffc0350fb8 pa=80350fb8 page_pa=80350000 | first_page_pa=8034f000
[SLUB] after free first: partial=2 full=0
[SLUB] before count(partial) after free
[SLUB] partial count=2
[SLUB] size=64 test done.
[SLUB] size=128 selected cache: objsize=128, slab_pages=1
[SLUB] new cache test for size=128 begin

```

size=32 size=64

多次“extra alloc”输出中：

```

extra alloc: va=ffffffffc0350fb8 pa=80350fb8 page_pa=80350000 |
first_page_pa=8034f000

```

可见 `page_pa` 由 8034f000 → 80350000，正好相差一个页 (0x1000)；新分配对象的物理页与旧 slab 相邻。

这符合 **buddy system** 连续页分配的预期行为，说明 SLUB 在分页 slab 与页级分配之间的接口协同是正确的。

```

[SLUB] partial count=2
[SLUB] size=64 test done.
[SLUB] size=128 selected cache: objsize=128, slab_pages=1
[SLUB] new cache test for size=128 begin
[SLUB] first alloc: va=ffffffffc0351f38 pa=80351f38 page_pa=80351000 | partial=1 full=0
[SLUB] capacity per slab = 31 objects
[SLUB] alloc i=15: va=ffffffffc03517b8 pa=803517b8 page_pa=80351000 | partial=1 full=0
[SLUB] alloc i=30: va=ffffffffc0351038 pa=80351038 page_pa=80351000 | partial=0 full=1
[SLUB] after fill: partial=0 full=1
[SLUB] extra alloc: va=ffffffffc0352f38 pa=80352f38 page_pa=80352000 | first_page_pa=80351000
[SLUB] after free first: partial=2 full=0
[SLUB] before count(partial) after free
[SLUB] partial count=2
[SLUB] size=128 test done.
[SLUB] size=256 selected cache: objsize=256, slab_pages=1
[SLUB] new cache test for size=256 begin
[SLUB] first alloc: va=ffffffffc0353e38 pa=80353e38 page_pa=80353000 | partial=1 full=0
[SLUB] capacity per slab = 15 objects
[SLUB] alloc i=7: va=ffffffffc0353738 pa=80353738 page_pa=80353000 | partial=1 full=0
[SLUB] alloc i=14: va=ffffffffc0353038 pa=80353038 page_pa=80353000 | partial=0 full=1
[SLUB] after fill: partial=0 full=1
[SLUB] extra alloc: va=ffffffffc0354e38 pa=80354e38 page_pa=80354000 | first_page_pa=80353000
[SLUB] after free first: partial=2 full=0
[SLUB] before count(partial) after free
[SLUB] partial count=2
[SLUB] size=256 test done.
[SLUB] size=512 selected cache: objsize=512, slab_pages=1
[SLUB] new cache test for size=512 begin
[SLUB] first alloc: va=ffffffffc0355c38 pa=80355c38 page_pa=80355000 | partial=1 full=0

```

size=128 size=256

也可以看到 partial / full 列表切换与地址对应正确

```
[SLUB] after fill: partial=0 full=1
```

```
[SLUB] after free first: partial=2 full=0
```


同时结合地址变化，在 slab 填满之前，所有对象的 `page_pa` 相同，`partial=1`；填满后，`partial` 变为 0，`full=1`；释放后又重新出现 `partial`。

这与 SLUB 算法的核心设计完全一致：**对象分配满 → 移入 full 列表；释放任意对象 → 移回 partial 列表。**

对应的地址连续性（相同 `page_pa`）和列表计数变化同步，说明链表操作与对象分配逻辑在页粒度上一致，**没有误判 slab 边界。**

```
[SLUB] partial count=2
[SLUB] size=256 test done.
[SLUB] size=512 selected cache: objsize=512, slab_pages=1
[SLUB] new cache test for size=512 begin
[SLUB] first alloc: va=ffffffffc0355c38 pa=80355c38 page_pa=80355000 | partial=1 full=0
[SLUB] capacity per slab = 7 objects
[SLUB] alloc i=3: va=ffffffffc0355638 pa=80355638 page_pa=80355000 | partial=1 full=0
[SLUB] alloc i=6: va=ffffffffc0355038 pa=80355038 page_pa=80355000 | partial=0 full=1
[SLUB] after fill: partial=0 full=1
[SLUB] extra alloc: va=ffffffffc0356c38 pa=80356c38 page_pa=80356000 | first_page_pa=80355000
[SLUB] after free first: partial=2 full=0
[SLUB] before count(partial) after free
[SLUB] partial count=2
[SLUB] size=512 test done.
[SLUB] size=1024 selected cache: objsize=1024, slab_pages=1
[SLUB] new cache test for size=1024 begin
[SLUB] first alloc: va=ffffffffc0357838 pa=80357838 page_pa=80357000 | partial=1 full=0
[SLUB] capacity per slab = 3 objects
[SLUB] alloc i=1: va=ffffffffc0357438 pa=80357438 page_pa=80357000 | partial=1 full=0
[SLUB] alloc i=2: va=ffffffffc0357038 pa=80357038 page_pa=80357000 | partial=0 full=1
[SLUB] after fill: partial=0 full=1
[SLUB] extra alloc: va=ffffffffc0358838 pa=80358838 page_pa=80358000 | first_page_pa=80357000
[SLUB] after free first: partial=2 full=0
[SLUB] before count(partial) after free
[SLUB] partial count=2
[SLUB] size=1024 test done.
[SLUB] size=2048 selected cache: objsize=2048, slab_pages=2
[SLUB] new cache test for size=2048 begin
```

size=512 size=1024

2. 双页 slab（2048B）验证

2048B 大小的对象使用双页 slab（`slab_pages=2`），测试输出显示：

```
[SLUB] size=1024 test done.
[SLUB] size=2048 selected cache: objsize=2048, slab_pages=2
[SLUB] new cache test for size=2048 begin
[SLUB] first alloc: va=ffffffffc035a038 pa=8035a038 page_pa=8035a000 | partial=1 full=0
[SLUB] capacity per slab = 3 objects
[SLUB] alloc i=1: va=ffffffffc0359838 pa=80359838 page_pa=80359000 | partial=1 full=0
[SLUB] alloc i=2: va=ffffffffc0359038 pa=80359038 page_pa=80359000 | partial=0 full=1
[SLUB] after fill: partial=0 full=1
[SLUB] extra alloc: va=ffffffffc035c038 pa=8035c038 page_pa=8035c000 | first_page_pa=8035a000
[SLUB] after free first: partial=1 full=1
[SLUB] before count(partial) after free
[SLUB] partial count=1
[SLUB] size=2048 test done.
[SLUB] size=4096 selected cache: objsize=4096, slab_pages=1
```

slab 初始化容量为 3 个对象；分配至满后，`partial=0 full=1`；释放首对象后，`partial` 和 `full` 同时存在（`partial=1 full=1`），说明 **多页 slab 的分配与状态切换逻辑稳定**，且未出现早期版本中的定位错误（`capacity=0`）问题。

3. 页级回退与特例检查

当对象大小恰好等于一页（4096B）时，程序走 **页级路径** 而非对象级 SLUB 缓存。该部分成功分配并释放整页，验证了 **页级回退机制** 的正确性。说明 SLUB 能在对象尺寸超出支持范围时正确退回底层 Buddy System。


```
[SLUB] size=4096 selected cache: objsize=4096, slab_pages=1
[SLUB] 4096B special case test begin
[SLUB] 4096B alloc: va=ffffffffc035d000 pa=8035d000 page_pa=8035d000
[SLUB] 4096B free done
```

4. 批量分配与复用验证

在 64B 对象的批量测试中：

一次性分配 1024 个对象后，出现 `partial=1 full=17`，说明系统自动创建了多个 slab 并管理良好；释放一半后 `partial=18 full=0`；再次分配（复用）后 `partial=1 full=17`，表明 SLUB 成功重用已释放对象，不发生碎片化，这一过程验证了 SLUB 的 **对象复用机制与空闲链表管理逻辑正确**。

```
[SLUB] 4096B free done
[SLUB] bulk alloc/free reuse test (size=64B)
[SLUB] arr[0]: va=ffffffffc034ffb8 pa=8034ffb8 page_pa=8034f000
[SLUB] arr[1]: va=ffffffffc0350f78 pa=80350f78 page_pa=80350000
[SLUB] arr[2]: va=ffffffffc0350f38 pa=80350f38 page_pa=80350000
[SLUB] arr[3]: va=ffffffffc0350ef8 pa=80350ef8 page_pa=80350000
[SLUB] arr[4]: va=ffffffffc0350eb8 pa=80350eb8 page_pa=80350000
[SLUB] after alloc: partial=1 full=17
[SLUB] after free half: partial=18 full=0
[SLUB] reuse[0]: va=ffffffffc036b078 pa=8036b078 page_pa=8036b000
[SLUB] reuse[1]: va=ffffffffc036b0f8 pa=8036b0f8 page_pa=8036b000
[SLUB] reuse[2]: va=ffffffffc036b178 pa=8036b178 page_pa=8036b000
[SLUB] after reuse: partial=1 full=17
[SLUB] big pages (2-page) fallback test begin
```

5. 大对象页级回退验证

对超过 4096B 的大对象测试显示，系统可以正确调用页级接口进行分配与释放。

```
[SLUB] reuse[2]: va=ffffffffc036b178 pa=8036b178 page_pa=8036b000
[SLUB] after reuse: partial=1 full=17
[SLUB] big pages (2-page) fallback test begin
[SLUB] big pages alloc: va=ffffffffc036d000 pa=8036d000
[SLUB] big pages free done
[SLUB] === slub_check() completed ===
check_alloc_page() succeeded!
satp virtual address: 0xffffffffc0207000
satp physical address: 0x0000000080207000
```

扩展练习Challenge：硬件的可用物理内存范围的获取方法（思考题）

- 如果 OS 无法提前知道当前硬件的可用物理内存范围，请问你有什么办法让 OS 获取可用物理内存范围？

方法一：通过固件接口获取内存布局

- 系统上电后，主板固件（如 BIOS 或 UEFI）会在自检 POST 阶段扫描物理内存，并维护一份精确的内存分布表，记录每个地址段的起始地址、长度和用途类型等。
- 当操作系统或引导程序启动时，可以调用主板固件提供的接口直接读取该内存地图。OS 只需解析这些条目，筛选出标记为可用的区域即可安全建立自己的物理页管理体系。

- 这种方法完全依赖固件，但缺点是必须在引导阶段调用，且依赖硬件厂商正确实现固件规范。

方法二：由引导加载器传递物理内存映射信息给内核

- 当操作系统本身无法直接访问固件接口，通常由引导加载器在实模式或早期阶段调用固件接口获取内存地图，并在加载内核时将该信息通过启动参数传递给内核。内核启动后直接解析这张表，即可获知哪些物理地址区间可安全使用。
- 这种方法实质上由引导程序完成底层硬件访问，内核只需信任并继承结果。其避免了内核与固件之间的复杂调用。

方法三：通过物理探测机制推测内存边界

- 当系统缺乏固件接口支持、也无引导器提供内存地图，操作系统可通过物理探测方式估算内存可用范围。其核心思路是从一个假定的起始地址开始，按页或大块递增访问物理地址，对每块内存执行读写测试，以判断其是否为真实的可读写内存。
- 若对某区域的写入能被正确读回且不会触发总线错误或机器检查异常，则推断该页属于有效 RAM；反之，则认为无效或保留区域。该方法能在完全缺乏系统描述表的环境下独立工作，但风险高，错误访问可能引发系统死机或损坏设备。

重要知识点总结

一、实验中的核心知识点及其与OS原理的对应关系

- 列出你认为本实验中重要的知识点，以及与对应的OS原理中的知识点，并简要说明你对二者的含义，关系，差异等方面的理解（也可能出现实验中的知识点没有对应的原理知识点）
- 本实验的重要知识点主要包括页Page结构管理、空闲内存链表的实现、First-Fit内存分配算法、空闲块合并机制等。这些内容在操作系统原理中分别对应于内存管理的分页机制与动态内存分配策略。
- 在操作系统原理中，内存分页是将逻辑内存划分为固定大小的页，并通过页表映射到物理页框，从而实现虚拟地址与物理地址的分离；而在本实验中，`struct Page` 结构体就是对物理页的抽象，每个 Page 对应一页物理内存，是物理内存管理的基本单元；实验中的 `free_list` 双向链表用于维护所有空闲页块的信息，这与操作系统中空闲区表的概念相对应，只是实现上采用了链表结构以支持动态插入和删除；至于 First-Fit 算法，其思想是“从头开始查找第一个能满足请求的空闲块”。实验中通过 `default_alloc_pages()` 的链表遍历来实现；实验中实现了 `default_free_pages()` 函数，对应 OS 中内存回收与碎片整理的过程，保证了内存释放与块合并。

- 第1个challenge (buddy system) 最重要的知识点当然就是buddy system算法，这部分知识点可以直接对应到PPT `4.20S_MemoryManagement-分页机制-2025` 中的内容。PPT中的内容重理论，学习完其中的内容可以对整个算法有直观的理解。实验中的内容重实践，在用代码实现buddy system的过程中可以加深对算法的理解，关注到一些细节。除此之外，我在实现时参考了wuwenbin的buddy system极简实现，其中的一些巧妙的设计给我留下了比较深刻的印象（**比如使用加权，即longest，来标记内存空闲状态，而不是一般的有限状态机**），这是OS原理课上不会涉及到的。
- 第2个challenge (slub) 上课时没有过多提及，所以在下一部分体现。

二、操作系统原理中重要但实验未涉及的知识点

- 列出你认为OS原理中很重要，但在实验中没有对应上的知识点

1. 虚拟内存管理机制 (Virtual Memory Management)

即通过页表 (Page Table)、TLB (Translation Lookaside Buffer) 等结构实现按需分配与地址空间隔离。这一机制是现代操作系统的核心，但实验中的 physical memory manager 仅在物理层面进行页分配，并未涉及虚拟地址到物理地址的映射。

2. 交换 (Swapping) 与页面置换算法 (Page Replacement)

例如 FIFO、LRU、Clock 等策略，这些算法用于在物理内存不足时将部分页换出到外存。

3. SLUB 内存分配机制

SLUB (**The Unqueued Slab Allocator**) 是 Linux 内核中用于高效管理小对象内存的分配器。它在伙伴系统 (Buddy System) 之上建立了对象级的分配层，将一页或多页内存切分为多个固定大小的对象块，用于频繁创建与销毁的小对象。

其核心思想是：

- **分层设计**：页级内存由 Buddy System 管理，SLUB 负责在页内细粒度地管理对象。
- **缓存复用**：为常见对象大小（如 32B、64B、128B 等）建立对应的缓存（`kmem_cache`），减少重复的初始化和碎片化。
- **快速分配与回收**：通过维护每个缓存的部分使用（`slabs_partial`）、已满（`slabs_full`）等链表，SLUB 能在常数时间内完成对象的分配与释放。

相比传统的 SLAB 算法，SLUB 省去了队列和引用计数等结构，简化了路径，提高了分配速度与并行性能。