

# ChCore 实验 2

班级：F2103003 姓名：黄伊新 学号：521030910070

## 1. 完成 buddy.c 中的 split\_chunk、merge\_chunk、buddy\_get\_pages、和 buddy\_free\_pages 函数

首先是 split\_chunk 函数，此处使用递归方法来进行实现。base case 是当传入的块的阶和 order 相同时，则直接返回该块，完成分裂。若不等，则需要将该块一分为二，并将其中的一块加入到 free list 中。在将该 chunk 的阶减 1 后即可通过 get\_buddy\_chunk 得到分裂后的后面一半的伙伴块，将其加入 free list 并进行有关设置。

```
1 static struct page *split_chunk(struct phys_mem_pool *pool, int
   order, struct page *chunk)
2 {
3     struct list_head *free_list;
4     struct page *buddy_page;
5     if (chunk->order == order){ // if this is the wanted order
6         return chunk;
7     }
8     chunk->order -= 1;
9
10    free_list = &(pool->free_lists[chunk->order].free_list);
11    buddy_page = get_buddy_chunk(pool, chunk); // split and get the
        remainder
12    buddy_page->order = chunk->order;
13    buddy_page->allocated = 0;
14    list_append(&(buddy_page->node), free_list); // add the
        buddy_page to the free list
15    pool->free_lists[chunk->order].nr_free += 1; // add one
16    return split_chunk(pool, order, chunk); // recursively split the
        chunk
17 }
18
```

第二个实现的时 merge\_chunk 函数。首先获得其对应的 buddy chunk，若传入 chunk 的阶没有超过最大限制，且能够找到一个空闲的伙伴块且阶数相同，则进行合并，若否则直接进行返回。需要注意的是在 merge 结束后是需要以编号更小的块作为新块的标识，这样才能通过得到块的位置来正确使用这个块的空间。

```
1 static struct page *merge_chunk(struct phys_mem_pool *pool, struct
   page *chunk)
2 {
3     struct page *buddy;
```

```

4 buddy = get_buddy_chunk(pool, chunk);
5 if (chunk->order == BUDDY_MAX_ORDER - 1 || buddy == NULL ||
    buddy->allocated == 1 || buddy->order != chunk->order){
6     return chunk;
7 }
8 list_del(&(buddy->node));
9 pool->free_lists[buddy->order].nr_free -= 1;
10 if ((u64) chunk < (u64) buddy){
11     buddy->allocated = 1;
12     chunk->order += 1; // already merge page with chunk
13     return merge_chunk(pool, chunk);
14 } else {
15     chunk->allocated = 1;
16     buddy->order += 1;
17     return merge_chunk(pool, buddy);
18 }
19 }
20

```

接下来完成的是 `buddy_get_page` 函数的实现。此处首先从 free list 中找到一个 best fit 的伙伴块，并进行分裂直到获得一个需要的阶数的伙伴块，并将该块标为已分配。

```

1 for (cur_order = order; cur_order < BUDDY_MAX_ORDER; ++cur_order)
2 {
3     free_list = &(pool->free_lists[cur_order].free_list);
4     if (pool->free_lists[cur_order].nr_free > 0) {
5         page = list_entry(pool->free_lists[cur_order].free_list.next,
6             struct page, node);
7         list_del(&(page->node));
8         pool->free_lists[cur_order].nr_free -= 1;
9         break;
10    }
11 }
12 if (page == NULL){
13     unlock(&pool->buddy_lock);
14     return NULL;
15 }
16 page = split_chunk(pool, order, page); // do the split
17 page->allocated = 1; // mark the page as allocated
18 page->order = order;

```

最后实现的是 `buddy_free_page` 的有关部分。这里主要需要处理一些特殊情况。若是正常的物理页且未分配则直接进行 merge，并将 merge 后的块加入 free list 并计数即可。

```

1 if (page == NULL) {return;}
2 if (page->allocated == 0) {return;}
3 page->allocated = 0; // marked as free

```

```

4 page = merge_chunk(pool, page); // merge chunk and return the
   merged one
5 order = page->order; // the order of the merged chunk
6 free_list = &(pool->free_lists[order].free_list);
7 list_append(&(page->node), free_list); // insert the free page
   into the free list
8 pool->free_lists[order].nr_free += 1; // do the updates
9

```

## 2. 练习题 2: 完成 kernel/mm/slab.c 中的 choose\_new\_current\_slab、alloc\_in\_slab\_impl 和 free\_in\_slab 函数

第一个实现的是 choose\_new\_current\_slab 函数。主要逻辑是若 partial list 中没有有空余的 slab 则从 buddy system 中直接获取一个新的，若有则取出并进行分配。从 partial list 中转换成 slab 使用的是 list\_entry 这个宏。

```

1 static void choose_new_current_slab(struct slab_pointer *pool, int
   order)
2 {
3     if (list_empty(&pool->partial_slab_list)) { // all slabs full
4         pool->current_slab = init_slab_cache(order, SIZE_OF_ONE_SLAB);
5         // get one from buddy system
6     } else {
7         pool->current_slab = list_entry(pool->partial_slab_list.next,
8         struct slab_header, node); // change current slab to the first
9         in the partial slab list
10        list_del(&pool->current_slab->node); // remove it from the
11        list
12    }
13 }
14

```

第二个是实现 alloc\_in\_slab\_impl 函数中的内容。此处首先判断目前使用的 slab 中是否有空闲，若否则从 partial list 中选择一个作为新的 slab。接着，从 current slab 中取出一个空闲的 slot，并进行有关更新。

```

1 if (current_slab->free_list_head == NULL){ // no free slot in
   current slab
2     choose_new_current_slab(&slab_pool[order], order); // choose a
   new slab
3     current_slab = slab_pool[order].current_slab; // change to the
   new slab
4 }
5 free_list = (struct slab_slot_list *)current_slab->free_list_head;
   // the first free slot in the current slab
6 next_slot = free_list->next_free; // the second free slot/NULL
7 current_slab->free_list_head = next_slot; // remove the slot from
   the slab
8 current_slab->current_free_cnt -= 1; // update the counts
9

```

最后是 free\_in\_slab 函数的有关实现。本处直接将 slot 插入回 free list，并计数即可。

```
1 slot->next_free = slab->free_list_head;
2 slab->free_list_head = (void *)slot; // insert slot to the
   free_list
3 slab->current_free_cnt += 1; // update the number of free slots
4
```

### 3. 练习题 3: 完成 \_\_kmalloc 函数，在适当位置调用对应函数，实现有关功能

本处主要是进行函数调用，若 size 较小则调用 slab 来进行分配，若较大则首先通过函数来计算出 best fit 的阶，然后使用 buddy system 进行分配。

```
1 if (size <= SLAB_MAX_SIZE) {
2     addr = alloc_in_slab(size, real_size);
3     #if ENABLE_MEMORY_USAGE_COLLECTING == ON
4         if(is_record && collecting_switch) {
5             record_mem_usage(*real_size, addr);
6         }
7     #endif
8 } else {
9     order = size_to_page_order(size);
10    addr = get_pages(order);
11 }
12
```

### 4. 练习题 4: 完成 query\_in\_pgtbl, map\_range\_in\_pgtbl\_common, 和 unmap\_range\_in\_pgtbl 和 mprotect\_in\_pgtbl 函数

首先是 query\_in\_pgtbl 函数。实现的思路大致是不断进行 get\_next\_ptp 函数的调用，若返回为-ENOMAPPING，即未映射，则直接进行返回；若碰到的是块描述符，则返回块描述符的内容和地址；若否则接着进行下一级页表的访问。本初因为代码较长就不放入 report 中了。

接着是 map\_range\_in\_pgtbl\_common 函数。该函数完成的任务是将一个区间内的虚拟地址空间映射到一个区间内的物理地址上。首先需要计算需要几个物理页来承载这些地址。接着依次来分配这些物理页。具体来说是通过 get\_next\_ptp 来获得/分配页表页，并具体来填写每一个页表项，并对已填写的页表项进行计数，若已达目标则退出。

然后是 unmap\_range\_in\_pgtbl 函数的填写。本处的进行方式同上面函数类似，都是先计算需要删除映射的数目，并依次进行。需要特别注意的是本处可能遇到本身就有没映射的情况，那么在每一层获取 page 的时候都需要额外进行判断，若该层本来就没有进行映射那么就不用再删除其下的项。若是最后发现有这个映射，那只需要在 l3 页表中将 valid 位设 0 就行。

最后是 mprotect\_in\_pgtbl 函数。这里主要是对最后一层的权限进行修改。方式和前面类似，本处直接到最后一层，并进行权限的修改即 set\_pte\_flags。

5. **思考题 5：阅读 Arm Architecture Reference Manual，思考要在操作系统中支持写时拷贝需要配置页表描述符的哪个字段，并在发生页错误时怎么处理**

在操作系统中支持写时拷贝需要配置页表描述符的 AP 字段，在 manual 的 B4.3 这一板块中。AP 即 Access Permission，它指定了访问该页的权限。在写时拷贝中，AP 字段应该设置为只读，以便多个进程可以共享同一块内存。

当进程试图写入只读页时，会触发一个页错误。当操作系统发现本处是由于本处是由于写时拷贝而触发的异常 Permission Fault 时，其会分配一个新的物理页，并将该页的内容复制到新的物理页中。然后，操作系统会更新页表，以便该页指向新的物理页。这样，该进程就可以修改其私有副本，而不会影响其他进程。

6. **思考题 6：在 ChCore 实验 Lab1 中没有为内核页表使用细粒度的映射，而是直接沿用了启动时的粗粒度页表，请思考这样做有什么问题？**

这导致了后续只能分配一个较大的页，较大的页可能会产生较多的内部碎片，内存的利用率较低。

7. **挑战题 7：使用前面实现的函数，配置细粒度页表**

挑战题本处没有在代码中加入这一部分，仅在报告中呈现一下大概的方式。

```
1 void reconfig_kernel_page_table() {
2     // memory
3     vmr_prop_t flags = AARCH64_MMU_ATTR_PAGE_AP_HIGH_RW_ELO_RW |
4     AARCH64_MMU_ATTR_PAGE_UX | NORMAL_MEMORY;
5     void *pgtbl;
6     pgtbl = get_pages(0);
7     memset(pgtbl, 0, PAGE_SIZE);
8     int ret;
9     size_t len1 = (1 << 30) - (1 << 24);
10    map_range_in_pgtbl(pgtbl, KBASE, 0, len1, flags);
11
12    // device memory
13    flags = AARCH64_MMU_ATTR_PAGE_AP_HIGH_RW_ELO_RW |
14    AARCH64_MMU_ATTR_PAGE_UX | DEVICE_MEMORY;
15    size_t len2 = 1 << 24;
16    map_range_in_pgtbl(pgtbl, KBASE + len1, len1, len2, flags);
17
18    // unmap
19    set_page_table(pgtbl);
20 }
```

8. **练习题 8：完成 do\_page\_fault 函数，将缺页异常进行转发**

这里直接调用 handle\_trans\_fault 函数即可，传入的函数就是本进程对应的 vmSPACE 和发生错误的地址。

```

1 ret = handle_trans_fault(current_thread->vmSPACE, fault_addr);
2

```

## 9. 练习题 9: 完成 find\_vmr\_for\_va 函数, 找到一个虚拟地址在其虚拟地址空间中的 VMR

这里是通过红黑树的搜索函数来获得包含该虚拟地址的 vmregion, 调用的函数为 rb\_search, 使用的是函数 cmp\_vmr\_and\_va 来进行比较。最后通过 rb\_entry 函数来获得包含该节点的 vmregion 即为我们所找的 vmregion。

```

1 struct vmregion *find_vmr_for_va(struct vmSPACE *vmSPACE, vaddr_t
  addr)
2 {
3     struct rb_root *vmr_root = &vmSPACE->vmr_tree;
4     struct rb_node *rb;
5     rb = rb_search(vmr_root, (const void*)addr, cmp_vmr_and_va);
6     if (rb == NULL) {
7         return NULL;
8     }
9     return rb_entry(rb, struct vmregion, tree_node);
10 }
11

```

## 10. 练习题 10: 完成 handle\_trans\_fault 函数, 实现按需物理页分配

若 pa==0, 则可以判断这一部分是需要进行按需分配的, 则首先通过 get\_pages 来获得一个新的物理页及其物理地址, 接着再将这一部分的物理页加入到页表中即可。

若否, 那么这里是被移出的页, 根据提示我们需要做的就是就是将这一部分页表重新进行映射, 调用 map\_range\_in\_ptbl 即可。

```

1 if (pa == 0) {
2     long rss = 0;
3     /* LAB 2 TODO 7 BEGIN */
4     /* BLANK BEGIN */
5     /* Hint: Allocate a physical page and clear it to 0. */
6     pa = virt_to_phys((vaddr_t)get_pages(0));
7     memset(phys_to_virt(pa), 0, PAGE_SIZE);
8     /* BLANK END */
9     /*
10      * Record the physical page in the radix tree:
11      * the offset is used as index in the radix tree
12      */
13     kdebug("commit: index: %ld, 0x%lx\n", index, pa);
14     commit_page_to_pmo(pmo, index, pa);
15
16     /* Add mapping in the page table */
17     lock(&vmSPACE->ptbl_lock);
18     /* BLANK BEGIN */
19     map_range_in_ptbl(vmSPACE->ptbl, fault_addr, pa, PAGE_SIZE,
    perm, rss);

```

```

20     /* BLANK END */
21     vmSPACE->rss += rss;
22     unlock(&vmSPACE->pgtbl_lock);
23 } else {
24     if (pmo->type == PMO_SHM || pmo->type == PMO_ANONYM) {
25         /* Add mapping in the page table */
26         long rss = 0;
27         lock(&vmSPACE->pgtbl_lock);
28         /* BLANK BEGIN */
29         map_range_in_pgtbl(vmSPACE->pgtbl, fault_addr, pa,
30             PAGE_SIZE, perm, rss);
31
32         /* BLANK END */
33         /* LAB 2 TODO 7 END */
34         vmSPACE->rss += rss;
35         unlock(&vmSPACE->pgtbl_lock);
36     }
37 }

```

## 11. 挑战题 11：简单实现换页技巧

本处选择最简单的方式实现换页。假设作为内存的物理页从 0 开始，而模拟硬盘的在高地址。本处认为共有 1G 作为内存的地址，页大小为 4K，故总共有  $2^{18}$  的页数。

接着使用数组来记录，即课本中提到的用一个位图来记录物理页是否有被映射，在换出时需要额外写一个函数来分配在指定区域内的一个物理页。具体实现如下。

```

1 #define page_num_in_memory (1 << 18) // number of pages in physical
   memory
2
3 bool is_mapped[page_num_in_memory];
4 vaddr_t origin_va[page_num_in_memory];
5 vmr_prop_t origin_perm[page_num_in_memory];
6 u64 access_time[page_num_in_memory]; // record the last time that
   it got accessed
7 u64 current_time = 0;
8
9
10 void alloc_in_memory(void *pgtbl, vaddr_t va, paddr_t pa,
   vmr_prop_t perm, long *rss) {
11     // a simple implementation of LRU
12     u64 min_time = -1;
13     u64 alloc_index;
14
15     for (u64 i = 0; i < page_num_in_memory; ++i) {
16         if (is_mapped[i]) {
17             if (min_time > access_time[i]) {
18                 min_time = access_time[i];

```

```

19         alloc_index = i;
20     }
21 }
22 else {
23     // empty page available
24     alloc_index = i;
25     break;
26 }
27 }
28
29 if (is_mapped[alloc_index]) {
30     // swap out
31     paddr_t pa_swap = get_pages_in_disk(); // to be
32 realized
33     unmap_range_in_pgtbl(pgtbl, origin_va[alloc_index],
34 PAGE_SIZE, rss);
35     memcpy(phys_to_virt(pa_swap), (void *)phys_to_virt
36 (alloc_index * PAGE_SIZE), PAGE_SIZE);
37     map_range_in_pgtbl(pgtbl, origin_va[alloc_index],
38 pa_swap, PAGE_SIZE, origin_perm[alloc_index], rss);
39 }
40
41 // allocate
42 is_mapped[alloc_index] = true;
43 origin_va[alloc_index] = va;
44 origin_perm[alloc_index] = perm;
45 access_time[alloc_index] = current_time++;
46
47 memcpy(phys_to_virt(alloc_index * PAGE_SIZE), (void *)
48 phys_to_virt(pa), PAGE_SIZE);
49 map_range_in_pgtbl(pgtbl, va, alloc_index * PAGE_SIZE,
50 PAGE_SIZE, perm, rss);
51 }
52

```