

内存管理

陈海波 / 夏虞斌

上海交通大学并行与分布式系统研究所

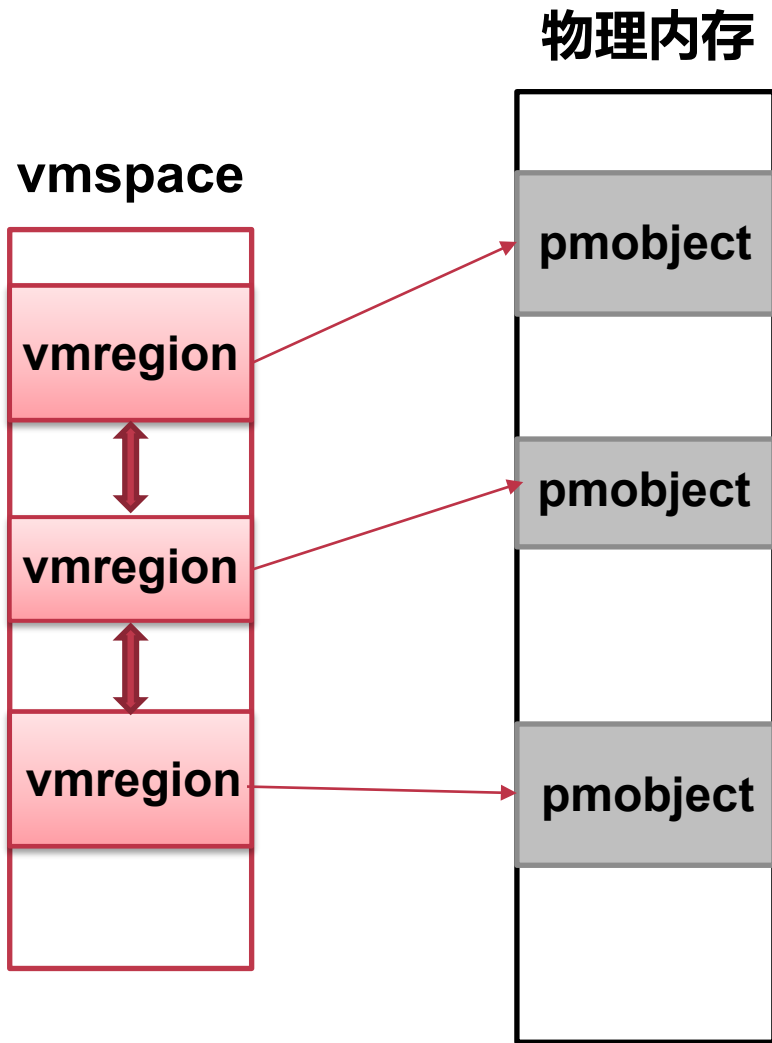
<https://ipads.se.sjtu.edu.cn>

版权声明

- 本内容版权归**上海交通大学并行与分布式系统研究所**所有
- 使用者可以将全部或部分本内容免费用于非商业用途
- 使用者在使用全部或部分本内容时请注明来源
 - 资料来自上海交通大学并行与分布式系统研究所+材料名字
- 对于不遵守此声明或者其他违法使用本内容者，将依法保留追究权
- 本内容的发布采用 Creative Commons Attribution 4.0 License
 - 完整文本：<https://creativecommons.org/licenses/by/4.0/legalcode>

核心数据结构

- 虚拟地址空间vmSPACE
- 虚拟地址区域vmregion
- 物理内存对象pobject



进程虚拟地址空间 vmSPACE (图)

- **vmr_list**
 - vmregion链表头
 - vmSPACE包含多个vmregion
- **pgtbl**
 - 页表基地址 (虚拟地址)
- **locks**
 - 多线程并发控制

```
struct vmSPACE {  
    /* List head of vmregion (vmr_list) */  
    struct list_head vmr_list;  
    /* Root page table */  
    vaddr_t *pgtbl;  
  
    /* The lock for manipulating vmregions */  
    struct lock vmSPACE_lock;  
    /* The lock for manipulating the page table */  
    struct lock pgtbl_lock;  
  
    /*  
     * For TLB flushing:  
     * Record the all the CPU that a vmSPACE ran on.  
     */  
    #ifdef CHCORE  
    u8 history_cpus[PLAT_CPU_NUM];  
    #endif  
  
    /* Heap-related: only used for user processes */  
    struct lock heap_lock;  
    struct vmregion *heap_vmr;  
    vaddr_t user_current_heap;  
  
    /* For the virtual address of mmap */  
    vaddr_t user_current_mmap_addr;  
};
```

进程虚拟地址区域vmregion

- node: 连接前后vmregion
- start: 区域起始位置
- size : 区域大小
- perm: 区域读写执行权限
- pmo: 对应的物理内存

```
struct vmregion {  
    struct list_head node; /* vmr_list */  
    vaddr_t start;  
    size_t size;  
    vmr_prop_t perm;  
    struct pobject *pmo;  
};
```

物理内存对象pobject

- 物理内存资源以对象的形式进行分配
 - 映射给vmregion

```
typedef u64 pmo_type_t;
#define PMO_ANONYM    0 /* lazy allocation */
#define PMO_DATA      1 /* immediate allocation */
#define PMO_FILE      2 /* file backed */
#define PMO_SHM       3 /* shared memory */
#define PMO_USER_PAGER 4 /* support user pager */
#define PMO_DEVICE    5 /* memory mapped device registers */

struct pobject {
    struct radix *radix; /* record physical pages */
    paddr_t start;
    size_t size;
    pmo_type_t type;
    atomic_cnt refcnt; /* free this pmo when refcnt is 0 */
};
```

物理内存对象pobject

- 每个类型的作用不同
 - PMO_ANONYM: 按需分配, 一般用于堆的分配
 - PMO_DATA: 立即分配, 比如用于加载代码和数据
 - PMO_FILE: 访问文件, 支持POSIX标准中的mmap
 - PMO_SHM: 共享内存
 - PMO_USER_PAGER: 支持用户态虚拟内存管理
 - PMO_DEVICE: 访问设备, 如MMIO和DMA

进程创建流程

- 首先创建虚拟地址空间vmospace
- 然后创建主线程
 - 分配栈空间（栈的vmregion和相应的pproject）
 - 加载代码和数据...

示例：为用户进程创建地址空间

- 创建并且初始化vmSPACE对象

```
274 vmSPACE = obj_alloc(TYPE_VMSPACE, sizeof(*vmSPACE));  
275 ▼ if (!vmSPACE) {  
276     r = -ENOMEM;  
277     goto out_free_obj_vmspace;  
278 }  
279 vmSPACE_init(vmSPACE);
```

初始化vmospace

- 初始化vmregion的链表
- 分配顶层页表页
- 体系结构相关的初始化
 - 如aarch64上可设置ASID
- 其它域的初始化

```
439 int vmospace_init(struct vmospace *vmospace)
440 {
441     init_list_head(&vmospace->vmr_list);
442     /* Allocate the root page table page */
443     vmospace->pgtbl = get_pages(0);
444     BUG_ON(vmspace->pgtbl == NULL);
445     memset((void*)vmospace->pgtbl, 0, PAGE_SIZE);
446
447     /* Architecture-dependent initialization */
448     arch_vmspace_init(vmspace);
449
450     /*
451      * Note: acquire vmospace_lock before pgtbl_lock
452      * when locking them together.
453      */
454     lock_init(&vmospace->vmospace_lock);
455     lock_init(&vmospace->pgtbl_lock);
456
457     vmospace->user_current_heap = HEAP_START;
458     lock_init(&vmospace->heap_lock);
459
460     /* The vmospace does not run on any CPU for now */
461     reset_history_cpus(vmspace);
462
463     /* Set the mmap area: this variable is protected
464      vmospace->user_current_mmap_addr = MMAP_START;
465
466     return 0;
467 }
```

示例：为第一个用户线程创建栈

- 分配并初始化物理内存对象，然后映射到地址空间

```
241     init_vmspace = obj_get(cap_group, VMSPACE_OBJ_ID, TYPE_VMSPACE);
242     obj_put(init_vmspace);
243
244     /* Allocate and setup a user stack for the init thread */
245     stack_pmo = obj_alloc(TYPE_PMO, sizeof(*stack_pmo));
246     if (!stack_pmo) {
247         ret = -ENOMEM;
248         goto out_fail;
249     }
250     pmo_init(stack_pmo, PMO_DATA, stack_size, 0);
251     stack_pmo_cap = cap_alloc(cap_group, stack_pmo, 0);
252     if (stack_pmo_cap < 0) {
253         ret = stack_pmo_cap;
254         goto out_free_obj_pmo;
255     }
256     ret = vmspace_map_range(init_vmspace, stack_base, stack_size,
257                             VMR_READ | VMR_WRITE, stack_pmo);
258     BUG_ON(ret != 0);
```

初始化pmobject

- 设置类型和大小
- 根据类型执行具体初始化:
 - 对于PMO_DATA类型, 直接分配物理内存
 - 对于PMO_ANONYM和PMO_SHM类型, 不会立即分配物理内存, 而只是创建radix_tree。将来按需分配的物理内存将被记录在这个树中。

```
498 /*
499  * Initialize an allocated pmobject.
500  * @paddr is only used when @type == PMO_DEVICE.
501  */
502 void pmobj_init(struct pmobject *pmo, pmobj_type_t type,
503                size_t len, paddr_t paddr)
504 {
505     memset((void*)pmo, 0, sizeof(*pmo));
506
507     len = ROUND_UP(len, PAGE_SIZE);
508     pmo->size = len;
509     pmo->type = type;
510
511     switch (type) {
512     case PMO_DATA: {
513         /*
514          * For PMO_DATA, the user will use it soon (we expect).
515          * So, we directly allocate the physical memory.
516          * Note that kmalloc(>2048) returns continous physical pages.
517          */
518         pmo->start = (paddr_t)virt_to_phys(kmalloc(len));
519         break;
520     }
521     case PMO_ANONYM:
522     case PMO_SHM: {
523         /*
524          * For PMO_ANONYM (e.g., stack and heap) or PMO_SHM,
525          * we do not allocate the physical memory at once.
526          */
527         pmo->radix = new_radix();
528         init_radix(pmo->radix);
529         break;
530     }
```

将pmbobject映射到地址空间中

- **vmospace_map_range**

1. 分配vmregion
2. 初始化vmregion
3. 将vmregion加入vmospace

```
167 int vmospace_map_range(struct vmospace *vmospace, vaddr_t va, size_t len,
168                        vmr_prop_t flags, struct pmbobject *pmo)
169 {
170     struct vmregion *vmr;
171     int ret;
172
173     /* skip other code */
174
175     vmr = alloc_vmrregion();
176     if (!vmr) {
177         ret = -ENOMEM;
178         goto out_fail;
179     }
180     vmr->start = va;
181     vmr->size = len;
182     vmr->perm = flags;
183     if (unlikely(pmo->type == PMO_DEVICE))
184         vmr->perm |= VMR_DEVICE;
185
186     /* Currently, one vmr has exactly one pmo */
187     vmr->pmo = pmo;
188
189     /*
190      * Note that each operation on the vmospace should be protected by
191      * the per_vmspace lock, i.e., vmospace_lock
192      */
193     lock(&vmospace->vmospace_lock);
194     ret = add_vmr_to_vmspace(vmspace, vmr);
195     unlock(&vmospace->vmospace_lock);
196     ...

```

将vmregion加入到vmospace中

- 检查vmregion是否合法，然后加入链表

```
78  /*
79  * This function should be surrounded with a lock.
80  *
81  * Modify the list with lock protection.
82  * Otherwise, concurrent operations may lead to erros.
83  */
84  static int add_vmr_to_vmospace(struct vmospace *vmospace, struct vmregion *vmr)
85  {
86      if (check_vmr_intersect(vmospace, vmr) != 0) {
87          kwarn("Detecting: vmr overlap\n");
88          BUG_ON(1);
89          return -EINVAL;
90      }
91
92      list_add(&(vmr->node), &(vmospace->vmr_list));
93      return 0;
94  }
```

Cont. 将pobject映射到地址空间中

- **vmSPACE_map_range**

1. 分配vmregion
2. 初始化vmregion
3. 将vmregion加入vmSPACE
4. 当pobject已经分配好物理内存了, 那么就直填写页表

```
203 int vmSPACE_map_range(struct vmSPACE *vmSPACE, vaddr_t va, size_t len,
204                       vmr_prop_t flags, struct pobject *pmo)
205 {
206     ... /* Omit some code */
207
208     /*
209      * Case-1:
210      * If the pmo type is PMO_DATA or PMO_DEVICE, we directly add mappings
211      * in the page table because the corresponding physical pages are
212      * prepared. In this case, early mapping avoids page faults and brings
213      * better performance.
214      *
215      * Case-2:
216      * Otherwise (for PMO_ANONYM and PMO_SHM), we use on-demand mapping.
217      * In this case, lazy mapping reduces the usage of physical memory resource.
218      */
219     if ((pmo->type == PMO_DATA) || (pmo->type == PMO_DEVICE))
220         fill_page_table(vmSPACE, vmr);
```

缺页异常

- **基于缺页异常，实现按需分配**
 - 例如，访问PMO_ANONYM映射的区域会触发缺页
- **(64位) 用户态程序发生缺页时下陷到内核**
 - 在x86_64上，触发 #PF （异常号13）
 - 在aarch64上，触发64位EL0同步异常 （异常号8）
 - 同时错误状态寄存器（ESR）中的错误代码（EC）是 ESR_EL1_EC_DABT_LEL 或 ESR_EL1_EC_IABT_LEL

```
18 #define ESR_EL1_EC_IABT_LEL      0b100000 /* Instruction Abort from a lower Exception level) */
19 #define ESR_EL1_EC_DABT_LEL      0b100100 /* Data Abort from a lower Exception level */
```


缺页异常处理函数

- 获取触发缺页的地址
 - aarch64: far_el1
 - X86_64: CR2
- 根据ESR中错误状态码具体处理
 - 翻译错误（无映射）
 - 权限错误

```
static inline vaddr_t get_fault_addr()
{
    vaddr_t addr;
    asm volatile ("mrs %0, far_el1\n\t" : "=r" (addr));
    return addr;
}
```

```
19 void do_page_fault(u64 esr, u64 fault_ins_addr)
20 {
21     vaddr_t fault_addr;
22     int fsc; // fault status code
23
24     fault_addr = get_fault_addr();
25     fsc = GET_ESR_EL1_FSC(esr);
26     switch (fsc) {
27         case DFSC_TRANS_FAULT_L0:
28         case DFSC_TRANS_FAULT_L1:
29         case DFSC_TRANS_FAULT_L2:
30         case DFSC_TRANS_FAULT_L3: {
31         int ret;
32
33         ret = handle_trans_fault(current_thread->vmSPACE, fault_addr);
34         break;
35     }
36     case DFSC_PERM_FAULT_L1:
37     case DFSC_PERM_FAULT_L2:
38     case DFSC_PERM_FAULT_L3:
39         /* Support COW later */
40
41         kinfo("do_page_fault: faulting ip is 0x%lx,"
42             "faulting address is 0x%lx,"
43             "fsc is perm_fault (0b%b)\n",
44             fault_ins_addr, fault_addr, fsc);
45         BUG_ON(1);
46         break;
47     }
```

处理按需分配导致的翻译错误

- 找到出错地址属于的 vmregion
- 拿到vmregion关联的pmo

```
28 int handle_trans_fault(struct vmpace *vmpace, vaddr_t fault_addr)
29 {
30     struct vmregion *vmr;
31     struct pmoobject *pmo;
32     paddr_t pa;
33     u64 offset;
34     u64 index;
35
36     /*
37      * Grab lock here.
38      * Because two threads (in same process) on different cores
39      * may fault on the same page, so we need to prevent them
40      * from adding the same mapping twice.
41      */
42
43     lock(&vmpace->vmpace_lock);
44     vmr = find_vmr_for_va(vmpace, fault_addr);
45     if (vmr == NULL) {
46         kinfo("handle_trans_fault: no vmr find for va 0x%lx!\n",
47             fault_addr);
48         kprint_vmr(vmpace);
49         // TODO: kill the process
50         kwarn("TODO: kill such faulting process.\n");
51         return -ENOMAPPING;
52     }
53
54     pmo = vmr->pmo;
```

Cont. 处理按需分配导致的翻译错误

- 分配物理页
- 在radix_tree中记录物理页信息
- 在页表中添加映射

```
54     pmo = vmr->pmo;
55     switch (pmo->type) {
56     case PMO_ANONYM:
57     case PMO_SHM: {
58         /* Simplified code */
59
60         pa = virt_to_phys(get_pages(0));
61         BUG_ON(pa == 0);
62         /*
63          * Record the physical page in the radix tree:
64          * the offset is used as index in the radix tree
65          */
66         kdebug("commit: index: %ld, 0x%lx\n", index, pa);
67         commit_page_to_pmo(pmo, index, pa);
68
69         /* Add mapping in the page table */
70         lock(&vmospace->pgtbl_lock);
71         fault_addr = ROUND_DOWN(fault_addr, PAGE_SIZE);
72         map_range_in_pgtbl(vmspace, fault_addr, pa,
73                             PAGE_SIZE, vmr->perm);
74         unlock(&vmospace->pgtbl_lock);
75     }
```