

对32-bit ARM而言，高于760M (0x2f80,0000) 的physical memory都属于highmem，kernel不能直接access，需要做kmap处理。

由struct page *page获得virtual address

如果是non-highmem

```
1.  #define page_address(page) lowmem_page_address(page)
2.
3.  static __always_inline void *lowmem_page_address(const struct page *page)
4.  {
5.      return __va(PFN_PHYS(page_to_pfn(page)));
6.  }
```

#define PFN_PHYS(x) ((phys_addr_t)(x) << PAGE_SHIFT)

如果是highmem

```
1.  void *kmap(struct page *page)
2.  {
3.      might_sleep();
4.      if (!PageHighMem(page))
5.          return page_address(page);
6.      return kmap_high(page);
7.  }
8.  EXPORT_SYMBOL(kmap);
```

核心函数kmap_high()

```

1.  /**
2.   * kmap_high - map a highmem page into memory
3.   * @page: &struct page to map
4.   *
5.   * Returns the page's virtual memory address.
6.   *
7.   * We cannot call this from interrupts, as it may block.
8.   */
9.  void *kmap_high(struct page *page)
10. {
11.     unsigned long vaddr;
12.
13.     /*
14.      * For highmem pages, we can't trust "virtual" until
15.      * after we have the lock.
16.      */
17.     lock_kmap();
18.     vaddr = (unsigned long)page_address(page);
19.     if (!vaddr)
20.         vaddr = map_new_virtual(page);
21.     pkmap_count[PKMAP_NR(vaddr)]++;
22.     BUG_ON(pkmap_count[PKMAP_NR(vaddr)] < 2);
23.     unlock_kmap();
24.     return (void*) vaddr;
25. }

```

linux/mm/highmem.c

```

1.  /**
2.   * page_address - get the mapped virtual address of a page
3.   * @page: &struct page to get the virtual address of
4.   *
5.   * Returns the page's virtual address.
6.   */
7.  void *page_address(const struct page *page)
8.  {
9.      unsigned long flags;
10.     void *ret;
11.     struct page_address_slot *pas;
12.
13.     if (!PageHighMem(page))
14.         return lowmem_page_address(page);
15.
16.     pas = page_slot(page);
17.     ret = NULL;
18.     spin_lock_irqsave(&pas->lock, flags);
19.     if (!list_empty(&pas->lh)) {
20.         struct page_address_map *pam;
21.
22.         list_for_each_entry(pam, &pas->lh, list) {
23.             if (pam->page == page) {
24.                 ret = pam->virtual;
25.                 goto done;
26.             }
27.         }
28.     }
29. done:
30.     spin_unlock_irqrestore(&pas->lock, flags);
31.     return ret;
32. }
33.
34. EXPORT_SYMBOL(page_address);

```

line 16

pas = page_slot(page);

```

1.  /*
2.   * Hash table bucket
3.   */
4.  static struct page_address_slot {
5.      struct list_head lh;           /* List of page_address_maps */
6.      spinlock_t lock;              /* Protect this bucket's list */
7.  } ____cacheline_aligned_in_smp page_address_hhtable[1<<PA_HASH_ORDER];
8.
9.  static struct page_address_slot *page_slot(const struct page *page)
10. {
11.     return &page_address_hhtable[hash_ptr(page, PA_HASH_ORDER)];
12. }

```

```
#define PA_HASH_ORDER 7
```

page_slot()由struct page *page作为hash key在page_address_hhtable hash table中得到one entry。如果为NULL,即该page还未在highmem中建立过mapping,则map_new_virtual(page) (后面分析)。

如果已经建立过mapping,

```

1.     if (!list_empty(&pas->lh)) {
2.         struct page_address_map *pam;
3.
4.         list_for_each_entry(pam, &pas->lh, list) {
5.             if (pam->page == page) {
6.                 ret = pam->virtual;
7.                 goto done;
8.             }
9.         }
10.    }

```

则在已有的mapping中查找到正确的struct page_address_map,返回该page对应的virtual.

```

1.  /*
2.   * Describes one page->virtual association
3.   */
4.  struct page_address_map {
5.      struct page *page;
6.      void *virtual;
7.      struct list_head list;
8.  };

```

page是位于highmem的struct page，而virtual则是该page所mapping的virtual address。凡是不同struct page，但hash key冲突的都被链接在一个list上。

```

1.      pkmap_count[PKMAP_NR(vaddr)]++;
2.      BUG_ON(pkmap_count[PKMAP_NR(vaddr)] < 2);

```

```

1.  static int pkmap_count[LAST_PKMAP];

```

```

1.  #define LAST_PKMAP          PTRS_PER_PTE

```

```

1.  #define PTRS_PER_PTE          512

```

pkmap_count[512]中记录了the specific page (in highmem) 's reference count.

目前 (3.18.7) kernel好像最多允许512 pages in highmem被mapping到virtual address.

如果在512 entries中找到一个empty entry,则

```

1.      vaddr = PKMAP_ADDR(last_pkmap_nr);
2.      set_pte_at(&init_mm, vaddr,
3.                  &(pkmap_page_table[last_pkmap_nr]), mk_pte(page, kmap_prot));
4.
5.      pkmap_count[last_pkmap_nr] = 1;
6.      set_page_address(page, (void *)vaddr);

```

last_pkmap_nr是0 to 512的pkmap_count[]的index.

```
#define PKMAP_ADDR(nr)          (PKMAP_BASE + ((nr) << PAGE_SHIFT))
```

```
#define PKMAP_BASE              (PAGE_OFFSET - PMD_SIZE)
```

$\Rightarrow 0xC000,0000 - 1UL \ll 21 = 3G - 2M = 0xBFE0,0000$

在Documentation/arm/memory.txt中有如下说明

PAGE_OFFSET high_memory-1 Kernel direct-mapped RAM region.

This maps the platforms RAM, and typically
maps all platform RAM in a 1:1 relationship.

PKMAP_BASE PAGE_OFFSET-1 Permanent kernel mappings

One way of mapping HIGHMEM pages into kernel
space.

即from 0xBFE0,0000 to 0xC000,0000 , 被kernel用于mapping来自highmem的内容。

There is the following log in kernel boot log

Virtual kernel memory layout:

vector : 0xffff0000 - 0xffff1000 (4 kB)

fixmap : 0xffc00000 - 0xffe00000 (2048 kB)

vmalloc : 0xf0000000 - 0xff000000 (240 MB)

lowmem : 0xc0000000 - 0xef800000 (760 MB)

pkmap : 0xbfe00000 - 0xc0000000 (2 MB)

modules : 0xbf000000 - 0xbfe00000 (14 MB)

`.text : 0xc0008000 - 0xc05cde80 (5912 kB)`

`.init : 0xc05ce000 - 0xc0602000 (208 kB)`

`.data : 0xc0602000 - 0xc0638728 (218 kB)`

`.bss : 0xc0638728 - 0xc06a8af4 (449 kB)`

那么one page in highmem怎么确定它被mapping的virtual address呢？

```
#define PKMAP_ADDR(nr)      (PKMAP_BASE + ((nr) << PAGE_SHIFT))
```

nr为该page在pkmap_count[512]数组中的index

`(nr) << PAGE_SHIFT`

把该index转换到相对[page的offset

把该offset + PKMAP_BASE , 即得到virtual address。

从上面的code上看到一个限制，当前active的可被利用的highmem最多为2M。

$4096 (1 \text{ page size}) * 512 = 2M。$

```
set_pte_at(&init_mm, vaddr,  
           &(pkmap_page_table[last_pkmap_nr]), mk_pte(page, kmap_prot));
```

```
pte_t *pkmap_page_table;
```

```
1. static void __init kmap_init(void)
2. {
3.     #ifdef CONFIG_HIGHMEM
4.         pkmap_page_table = early_pte_alloc(pmd_off_k(PKMAP_BASE),
5.             PKMAP_BASE, _PAGE_KERNEL_TABLE);
6.
7.         fixmap_page_table = early_pte_alloc(pmd_off_k(FIXADDR_START),
8.             FIXADDR_START, _PAGE_KERNEL_TABLE);
9.     #endif
10. }
```

pkmap_page_table是用来mapping kmap的pte array。

mk_pte(page, kmap_prot)生成pte.由page可以获得physical address , 而由physical address + 该page的access properities = pte value.

而set_pte_at()则是把pte与virtual address关联。同时kmap属于init_mm地址空间。

如果2M的kmap空间use up,那么在map_new_virtual()中


```

1.      /*
2.      * Sleep for somebody else to unmap their entries
3.      */
4.      {
5.          DECLARE_WAITQUEUE(wait, current);
6.          wait_queue_head_t *pkmap_map_wait =
7.              get_pkmap_wait_queue_head(color);
8.
9.          __set_current_state(TASK_UNINTERRUPTIBLE);
10.         add_wait_queue(pkmap_map_wait, &wait);
11.         unlock_kmap();
12.         schedule();
13.         remove_wait_queue(pkmap_map_wait, &wait);
14.         lock_kmap();
15.
16.         /* Somebody else might have mapped it while we slept */
17.         if (page_address(page))
18.             return (unsigned long)page_address(page);
19.
20.         /* Re-start */
21.         goto start;
22.     }

```

调用kmap()的process只能睡眠，等着其他process调用kunmap()，从而再次查询在pkmap_count[512]中是否有empty entry。

总结：

在atomic context中不能调用kmap() / kunmap(),因为该调用会sleep。

判断某个address的空间是否来自kmap(highmem)，只要判断其是否在0xBFE0,0000和

0xC000,0000之间。

```
1. struct page *kmap_to_page(void *vaddr)
2. {
3.     unsigned long addr = (unsigned long)vaddr;
4.
5.     if (addr >= PKMAP_ADDR(0) && addr < PKMAP_ADDR(LAST_PKMAP)) {
6.         int i = PKMAP_NR(addr);
7.         return pte_page(pkmap_page_table[i]);
8.     }
9.
10.    return virt_to_page(addr);
11. }
12. EXPORT_SYMBOL(kmap_to_page);
```