# Memory Model -4- (APIs)

📅 2016-03-31 (http://jake.dothome.co.kr/mem_map_page/)    👤 Moon Young-il (http://jake.dothome.co.kr/author/admin/)    📁 Linux Kernel (http://jake.dothome.co.kr/category/linux/)

## Memory Model -4- (APIs)

### for_each_migratetype_order()

include/linux/mmzone.h

```
1  #define for_each_migratetype_order(order, type) \
2          for (order = 0; order < MAX_ORDER; order++) \
3                  for (type = 0; type < MIGRATE_TYPES; type++)
```

- Loop around the number of MAX_ORDER (11) used by the buddy memory allocator.
- Loop as much as you MIGRATE_TYPES responsible for managing migration flags for memory hotplugs.

## Zones and Nodes

### is_highmem_idx()

include/linux/mmzone.h

```
1  static inline int is_highmem_idx(enum zone_type idx)
2  {
3  #ifdef CONFIG_HIGHMEM
4          return (idx == ZONE_HIGHMEM ||
5                  (idx == ZONE_MOVABLE && movable_zone == ZONE_HIGHMEM));
6  #else
7          return 0;
8  #endif
9  }
```

### zone_idx()

include/linux/mmzone.h

```
1  /*
2   * zone_idx() returns 0 for the ZONE_DMA zone, 1 for the ZONE_NORMAL zone, etc.
3   */
4  #define zone_idx(zone)          ((zone) - (zone)->zone_pgdat->node_zones)
```

- Returns the zone index number.

- e.g. 0 and 1 are returned when using ZONE_DMA and ZONE_NORMAL.
- e.g. 0 is returned if only ZONE_NORMAL is used.

## set_page_links()

include/linux/mm.h

```
1  static inline void set_page_links(struct page *page, enum zone_type zon
   e,
2          unsigned long node, unsigned long pfn)
3  {
4          set_page_zone(page, zone);
5          set_page_node(page, node);
6  #ifdef SECTION_IN_PAGE_FLAGS
7          set_page_section(page, pfn_to_section_nr(pfn));
8  #endif
9  }
```

- Set the zone, node, and section information in page->flags.

## set_page_zone()

include/linux/mm.h

```
1  static inline void set_page_zone(struct page *page, enum zone_type zone)
2  {
3          page->flags &= ~(ZONES_MASK << ZONES_PGSHIFT);
4          page->flags |= (zone & ZONES_MASK) << ZONES_PGSHIFT;
5  }
```

- Set the zone information in page->flags.

## set_page_node()

include/linux/mm.h

```
1  static inline void set_page_node(struct page *page, unsigned long node)
2  {
3          page->flags &= ~(NODES_MASK << NODES_PGSHIFT);
4          page->flags |= (node & NODES_MASK) << NODES_PGSHIFT;
5  }
```

- page->flags.

## page_zone_id()

include/linux/mm.h

```
01  /*
02   * The identification function is mainly used by the buddy allocator for
03   * determining if two pages could be buddies. We are not really identify
    ing
04   * the zone since we could be using the section number id if we do not h
    ave
05   * node id available in page flags.
06   * We only guarantee that it will return the same value for two combinab
    le
```

```
07    * pages in a zone.
08    */
09   static inline int page_zone_id(struct page *page)
10   {
11          return (page->flags >> ZONEID_PGSHIFT) & ZONEID_MASK;
12   }
```

Extracts the zone id from the page and returns it.

---

## About the Sprsemem section

### set_page_section()

include/linux/mm.h

```
1   #ifdef SECTION_IN_PAGE_FLAGS
2   static inline void set_page_section(struct page *page, unsigned long section)
3   {
4          page->flags &= ~(SECTIONS_MASK << SECTIONS_PGSHIFT);
5          page->flags |= (section & SECTIONS_MASK) << SECTIONS_PGSHIFT;
6   }
7   #endif
```

- page->flags.

### __pfn_to_section()

include/linux/mmzone.h

```
1   static inline struct mem_section *__pfn_to_section(unsigned long pfn)
2   {
3          return __nr_to_section(pfn_to_section_nr(pfn));
4   }
```

Returns information about the mem_section structure that corresponds to the pfn value.

- pfn_to_section_nr()
  - Get the section number by the pfn value.
- __nr_to_section()
  - Returns mem_section struct information with a section number.

### pfn_to_section_nr()

include/linux/mmzone.h

```
1   #define pfn_to_section_nr(pfn) ((pfn) >> PFN_SECTION_SHIFT)
```

- Returns the sparse index of the pfn.
  - e.g. Realview-PBX
    - Since the section size is 256M units (PFN_SECTION_SHIFT=16), the section number is from 0~15

## __nr_to_section()

include/linux/mmzone.h

```
1  static inline struct mem_section *__nr_to_section(unsigned long nr)
2  {
3          if (!mem_section[SECTION_NR_TO_ROOT(nr)])
4                  return NULL;
5          return &mem_section[SECTION_NR_TO_ROOT(nr)][nr & SECTION_ROOT_MA
SK];
6  }
```

- Returns mem_section struct information with a section number.

## SECTION_NR_TO_ROOT()

include/linux/mmzone.h

```
1  #define SECTION_NR_TO_ROOT(sec) ((sec) / SECTIONS_PER_ROOT)
```

- Returns the ROOT number as the section number.

```
1  #ifdef CONFIG_SPARSEMEM_EXTREME
2  #define SECTIONS_PER_ROOT       (PAGE_SIZE / sizeof (struct mem_sectio
n))
3  #else
4  #define SECTIONS_PER_ROOT       1
5  #endif
```

- Number of sections per ROOT
  - PAGE_SIZE (4K) can fit mem_section structure

## present_section_nr()

include/linux/mmzone.h

```
1  static inline int present_section_nr(unsigned long nr)
2  {
3          return present_section(__nr_to_section(nr));
4  }
```

Make sure you have a mem_section that corresponds to the section number. If it is not ready, the section means hole.

- __nr_to_section()
  - Retrieves the mem_section struct information by section number.
- present_section()
  - mem_section Check if a section exists in the struct information.

## present_section()

include/linux/mmzone.h

```
1  static inline int present_section(struct mem_section *section)
```

```
2  {
3          return (section && (section->section_mem_map & SECTION_MARKED_PR
   ESENT));
4  }
```

mem_section Check if a section exists in the struct information.

- Make sure the SECTION_MARKED_PRESENT identification bit is set.

### __section_mem_map_addr()

include/linux/mmzone.h

```
1  static inline struct page *__section_mem_map_addr(struct mem_section *se
   ction)
2  {
3          unsigned long map = section->section_mem_map;
4          map &= SECTION_MAP_MASK;
5          return (struct page *)map;
6  }
```

Returns the mem_map address for the Sparse memory section.

include/linux/mmzone.h

```
01  /*
02   * We use the lower bits of the mem_map pointer to store
03   * a little bit of information.  There should be at least
04   * 3 bits here due to 32-bit alignment.
05   */
06  #define SECTION_MARKED_PRESENT  (1UL<<0)
07  #define SECTION_HAS_MEM_MAP     (1UL<<1)
08  #define SECTION_MAP_LAST_BIT    (1UL<<2)
09  #define SECTION_MAP_MASK        (~(SECTION_MAP_LAST_BIT-1))
10  #define SECTION_NID_SHIFT       2
```

# Enabling and disabling references on a page

### get_page()

include/linux/mm.h

```
01  static inline void get_page(struct page *page)
02  {
03          page = compound_head(page);
04          /*
05           * Getting a normal page or the head of a compound page
06           * requires to already have an elevated page->_refcount.
07           */
08          VM_BUG_ON_PAGE(page_ref_count(page) <= 0, page);
09          page_ref_inc(page);
10  }
```

Increments the reference counter by 1.

### get_page_unless_zero()

include/linux/mm.h

```
1   /*
2    * Try to grab a ref unless the page has a refcount of zero, return fals
    e if
3    * that is the case.
4    * This can be called when MMU is off so it must not access
5    * any of the virtual mappings.
6    */
```

```
1   static inline int get_page_unless_zero(struct page *page)
2   {
3           return page_ref_add_unless(page, 1, 0);
4   }
```

After reading the reference counter (p->_refcount), increment only if it differs from the 0 value. If the result value is not zero, it returns true.

## put_page()

include/linux/mm.h

```
01  static inline void put_page(struct page *page)
02  {
03          page = compound_head(page);
04
05          /*
06           * For devmap managed pages we need to catch refcount transition
    from
07           * 2 to 1, when refcount reach one it means the page is free and
    we
08           * need to inform the device driver through callback. See
09           * include/linux/memremap.h and HMM for details.
10           */
11          if (put_devmap_managed_page(page))
12                  return;
13
14          if (put_page_testzero(page))
15                  __put_page(page);
16  }
```

Reduces the reference counter by 1. If it reaches 0, the page will be retrieved.

## put_page_testzero()

include/linux/mm.h

```
01  /*
02   * Methods to modify the page usage count.
03   *
04   * What counts for a page usage:
05   * - cache mapping   (page->mapping)
06   * - private data    (page->private)
07   * - page mapped in a task's page tables, each mapping
08   *   is counted separately
09   *
10   * Also, many kernel routines increase the page count before a critical
11   * routine so they can be sure the page doesn't go away from under them.
12   */
13
14  /*
```

```
15    * Drop a ref, return true if the refcount fell to zero (the page has no
      users)
16    */

 1  static inline int put_page_testzero(struct page *page)
 2  {
 3          VM_BUG_ON_PAGE(page_ref_count(page) == 0, page);
 4          return page_ref_dec_and_test(page);
 5  }
```

Decrement the reference counter on the page and check if it is 0 (used complete) to return whether it is used or not.

- 0=in use, 1=in use (when the reference _count becomes 0)

## Page vs PFN Conversion

Conversion between PFN and page structure pointers uses two APIs:

include/asm-generic/memory_model.h

```
1  #define page_to_pfn __page_to_pfn
2  #define pfn_to_page __pfn_to_page
```

- page_to_pfn()
    - Page struct pointer to get the pfn value.
- pfn_to_page()
    - pfn value to get the pointer to the page structure.

As follows, the conversion method varies depending on the flat and sparse physics memory models, and in the case of the sparse physics model, it is further divided into two types depending on whether vmemmap is used or not.

### CONFIG_FLATMEM

```
1  #define __pfn_to_page(pfn)        (mem_map + ((pfn) - ARCH_PFN_OFFSET))
2  #define __page_to_pfn(page)       ((unsigned long)((page) - mem_map) + \
3                                     ARCH_PFN_OFFSET)
```

- __pfn_to_page()
    - ARCH_PFN_OFFSET refers to the starting PFN value of the physical DRAM.
    - mem_map[@pfn – Physics DRAM Start PFN]

### CONFIG_SPARSEMEM

```
01  /*
02   * Note: section's mem_map is encoded to reflect its start_pfn.
03   * section[i].section_mem_map == mem_map's address - start_pfn;
04   */
05  #define __page_to_pfn(pg)                                            \
06  ({      const struct page *__pg = (pg);                              \
07          int __sec = page_to_section(__pg);                           \
```

```
08          (unsigned long)(__pg - __section_mem_map_addr(__nr_to_section(__
   sec))); \
09  })
10
11  #define __pfn_to_page(pfn)                              \
12  ({      unsigned long __pfn = (pfn);                    \
13          struct mem_section *__sec = __pfn_to_section(__pfn);    \
14          __section_mem_map_addr(__sec) + __pfn;          \
15  })
```

- __pfn_to_page()
    - Replace pfn with sections, and then access mem_section[][] to return the mam_map[@pfn]
      address for the section.

### CONFIG_SPARSEMEM & CONFIG_SPARSEMEM_VMEMMAP

```
1  /* memmap is virtually contiguous.   */
2  #define __pfn_to_page(pfn)      (vmemmap + (pfn))
3  #define __page_to_pfn(page)     (unsigned long)((page) - vmemmap)
```

- __pfn_to_page()
    - = mem_map[@pfn]
    - vmemmap = mem_map[0].

# Page Flags

include/linux/page-flags.h

```
01  /*
02   * Various page->flags bits:
03   *
04   * PG_reserved is set for special pages. The "struct page" of such a pag
   e
05   * should in general not be touched (e.g. set dirty) except by its owne
   r.
06   * Pages marked as PG_reserved include:
07   * - Pages part of the kernel image (including vDSO) and similar (e.g. B
   IOS,
08   *   initrd, HW tables)
09   * - Pages reserved or allocated early during boot (before the page allo
   cator
10   *   was initialized). This includes (depending on the architecture) the
11   *   initial vmemmap, initial page tables, crashkernel, elfcorehdr, and
   much
12   *   much more. Once (if ever) freed, PG_reserved is cleared and they wi
   ll
13   *   be given to the page allocator.
14   * - Pages falling into physical memory gaps - not IORESOURCE_SYSRAM. Tr
   ying
15   *   to read/write these pages might end badly. Don't touch!
16   * - The zero page(s)
17   * - Pages not added to the page allocator when onlining a section becau
   se
18   *   they were excluded via the online_page_callback() or because they a
   re
19   *   PG_hwpoison.
20   * - Pages allocated in the context of kexec/kdump (loaded kernel image,
21   *   control pages, vmcoreinfo)
```

```
22    * - MMIO/DMA pages. Some architectures don't allow to ioremap pages tha
      t are
23    *   not marked PG_reserved (as they might be in use by somebody else wh
      o does
24    *   not respect the caching strategy).
25    * - Pages part of an offline section (struct pages of offline sections
      should
26    *   not be trusted as they will be initialized when first onlined).
27    * - MCA pages on ia64
28    * - Pages holding CPU notes for POWER Firmware Assisted Dump
29    * - Device memory (e.g. PMEM, DAX, HMM)
30    * Some PG_reserved pages will be excluded from the hibernation image.
31    * PG_reserved does in general not hinder anybody from dumping or swappi
      ng
32    * and is no longer required for remap_pfn_range(). ioremap might requir
      e it.
33    * Consequently, PG_reserved for a page mapped into user space can indic
      ate
34    * the zero page, the vDSO, MMIO pages or device memory.
35    *
36    * The PG_private bitflag is set on pagecache pages if they contain file
      system
37    * specific data (which is normally at page->private). It can be used by
38    * private allocations for its own usage.
39    *
40    * During initiation of disk I/O, PG_locked is set. This bit is set befo
      re I/O
41    * and cleared when writeback _starts_ or when read _completes_. PG_writ
      eback
42    * is set before writeback starts and cleared when it finishes.
43    *
44    * PG_locked also pins a page in pagecache, and blocks truncation of the
      file
45    * while it is held.
46    *
47    * page_waitqueue(page) is a wait queue of all tasks waiting for the pag
      e
48    * to become unlocked.
49    *
50    * PG_swapbacked is set when a page uses swap as a backing storage.  Thi
      s are
51    * usually PageAnon or shmem pages but please note that even anonymous p
      ages
52    * might lose their PG_swapbacked flag when they simply can be dropped
      (e.g. as
53    * a result of MADV_FREE).
54    *
55    * PG_uptodate tells whether the page's contents is valid.  When a read
56    * completes, the page becomes uptodate, unless a disk I/O error happene
      d.
57    *
58    * PG_referenced, PG_reclaim are used for page reclaim for anonymous and
59    * file-backed pagecache (see mm/vmscan.c).
60    *
61    * PG_error is set to indicate that an I/O error occurred on this page.
62    *
63    * PG_arch_1 is an architecture specific page state bit.  The generic co
      de
64    * guarantees that this bit is cleared for a page when it first is enter
      ed into
65    * the page cache.
66    *
67    * PG_hwpoison indicates that a page got corrupted in hardware and conta
      ins
68    * data with incorrect ECC bits that triggered a machine check. Accessin
      g is
69    * not safe since it may cause another machine check. Don't touch!
70    */
```

```
01  /*
02   * Don't use the pageflags directly.  Use the PageFoo macros.
03   *
04   * The page flags field is split into two parts, the main flags area
05   * which extends from the low bits upwards, and the fields area which
06   * extends from the high bits downwards.
07   *
08   *  | FIELD | ... | FLAGS |
09   *  N-1           ^       0
10   *             (NR_PAGEFLAGS)
11   *
12   * The fields area is reserved for fields mapping zone, node (for NUMA) and
13   * SPARSEMEM section (for variants of SPARSEMEM that require section ids like
14   * SPARSEMEM_EXTREME with !SPARSEMEM_VMEMMAP).
```

```
01  enum pageflags {
02          PG_locked,                  /* Page is locked. Don't touch. */
03          PG_referenced,
04          PG_uptodate,
05          PG_dirty,
06          PG_lru,
07          PG_active,
08          PG_workingset,
09          PG_waiters,                 /* Page has waiters, check its waitqueue. Must be bit #7 and in the same byte as "PG_locked" */
10          PG_error,
11          PG_slab,
12          PG_owner_priv_1,            /* Owner use. If pagecache, fs may use*/
13          PG_arch_1,
14          PG_reserved,
15          PG_private,                 /* If pagecache, has fs-private data */
16          PG_private_2,               /* If pagecache, has fs aux data */
17          PG_writeback,               /* Page is under writeback */
18          PG_head,                    /* A head page */
19          PG_mappedtodisk,            /* Has blocks allocated on-disk */
20          PG_reclaim,                 /* To be reclaimed asap */
21          PG_swapbacked,              /* Page is backed by RAM/swap */
22          PG_unevictable,             /* Page is "unevictable"  */
23  #ifdef CONFIG_MMU
24          PG_mlocked,                 /* Page is vma mlocked */
25  #endif
26  #ifdef CONFIG_ARCH_USES_PG_UNCACHED
27          PG_uncached,                /* Page has been mapped as uncached */
28  #endif
29  #ifdef CONFIG_MEMORY_FAILURE
30          PG_hwpoison,                /* hardware poisoned page. Don't touch */
31  #endif
32  #if defined(CONFIG_PAGE_IDLE_FLAG) && defined(CONFIG_64BIT)
33          PG_young,
34          PG_idle,
35  #endif
36  #ifdef CONFIG_64BIT
37          PG_arch_2,
38  #endif
39  #ifdef CONFIG_KASAN_HW_TAGS
40          PG_skip_kasan_poison,
41  #endif
42          __NR_PAGEFLAGS,
43
44          /* Filesystems */
45          PG_checked = PG_owner_priv_1,
46
47          /* SwapBacked */
48          PG_swapcache = PG_owner_priv_1, /* Swap page: swp_entry_t in private */
```

```
49
50          /* Two page bits are conscripted by FS-Cache to maintain local c
   aching
51           * state.   These bits are set on pages belonging to the netfs's
   inodes
52           * when those inodes are being locally cached.
53           */
54          PG_fscache = PG_private_2,      /* page backed by cache */
55
56          /* XEN */
57          /* Pinned in Xen as a read-only pagetable page. */
58          PG_pinned = PG_owner_priv_1,
59          /* Pinned as part of domain save (see xen_mm_pin_all()). */
60          PG_savepinned = PG_dirty,
61          /* Has a grant mapping of another (foreign) domain's page. */
62          PG_foreign = PG_owner_priv_1,
63          /* Remapped by swiotlb-xen. */
64          PG_xen_remapped = PG_owner_priv_1,
65
66          /* SLOB */
67          PG_slob_free = PG_private,
68
69          /* Compound pages. Stored in first tail page's flags */
70          PG_double_map = PG_workingset,
71
72 #ifdef CONFIG_MEMORY_FAILURE
73          /*
74           * Compound pages. Stored in first tail page's flags.
75           * Indicates that at least one subpage is hwpoisoned in the
76           * THP.
77           */
78          PG_has_hwpoisoned = PG_mappedtodisk,
79 #endif
80
81          /* non-lru isolated movable page */
82          PG_isolated = PG_reclaim,
83
84          /* Only valid for buddy pages. Used to track pages that are repo
   rted */
85          PG_reported = PG_uptodate,
86 };
```

## Additional information logged in page->flags

linux/page-flags-layout.h

```
01 /*
02  * page->flags layout:
03  *
04  * There are five possibilities for how page->flags get laid out.  The f
   irst
05  * pair is for the normal case without sparsemem. The second pair is for
06  * sparsemem when there is plenty of space for node and section informat
   ion.
07  * The last is when there is insufficient space in page->flags and a sep
   arate
08  * lookup is necessary.
09  *
10  * No sparsemem or sparsemem vmemmap: |       NODE     | ZONE
   |             ... | FLAGS |
11  *      " plus space for last_cpuid: |       NODE     | ZONE | LAST_CPU
   PID ... | FLAGS |
12  * classic sparse with space for node:| SECTION | NODE | ZONE
   |         ... | FLAGS |
```

```
13    *        " plus space for last_cpupid: | SECTION | NODE | ZONE | LAST_CPU
      PID ... | FLAGS |
14    * classic sparse no space for node:  | SECTION |     ZONE    | ... | FL
      AGS |
15    */
```

Depending on the kernel configuration, page->flags records SECTION, NODE, ZONE, and LAST_CPUPID information in addition to flags.

## Reserved flag (example)

### PageReserved(), SetPageReserved(), ClearPageReserved(), __ClearPageReserved()

include/linux/page-flags.h

```
1  PAGEFLAG(Reserved, reserved) __CLEARPAGEFLAG(Reserved, reserved)
```

- PageReserved(), SetPageReserved(), ClearPageReserved(), and __ClearPageReserved() static inline functions are created.

```
1  #define PAGEFLAG(uname, lname) TESTPAGEFLAG(uname, lname)
   \
2          SETPAGEFLAG(uname, lname) CLEARPAGEFLAG(uname, lname)
```

- The following macros are used to create the PageXXX(), SetPageXXX(), and ClearPageXXX() static inline functions.

```
01  /*
02   * Macros to create function definitions for page flags
03   */
04  #define TESTPAGEFLAG(uname, lname)
    \
05  static inline int Page##uname(const struct page *page)
    \
06                      { return test_bit(PG_##lname, &page->flags); }
07
08  #define SETPAGEFLAG(uname, lname)
    \
09  static inline void SetPage##uname(struct page *page)
    \
10                      { set_bit(PG_##lname, &page->flags); }
11
12  #define CLEARPAGEFLAG(uname, lname)
    \
13  static inline void ClearPage##uname(struct page *page)
    \
14                      { clear_bit(PG_##lname, &page->flags); }
```

```
1  #define __CLEARPAGEFLAG(uname, lname)
   \
2  static inline void __ClearPage##uname(struct page *page)
   \
3                      { __clear_bit(PG_##lname, &page->flags); }
```

- test_bit()
  - &page->flags to find out whether the PG_xxxxx number bit has been set.
- set_bit()
  - &page->flags sets the PG_xxxxx number bit atomically.
- clear_bit()
  - &PG_xxxxx number bit of &page->flags is atomically cleared.
- __clear_bit()
  - Clear the PG_xxxxx number bit of &page->flags. (non-atomic)

## Regrouping of some flags (P->page_type)

The four PG_buddy, PG_ballon, PG_kmemcg, and PG_table flags below have been replaced with P->_mapcount, which is then declared as a union and used as a shared P->page_type.

- Flags that were initially managed by P->Flags were split into P->_mapcount.
- The new kernel then uses p->page_type shared as a union instead of p->_mapcount. However, since the initial value of p->_mapcount is -1 (0xffff_ffff), the bits are set and turned off in reverse.
  - 예) Set Buddy
    - Old kernel: p->_mapcount = PAGE_BUDDY_MAPCOUNT_VALUE(-128)
    - new 커널: p->page_type &= ~0x80
  - 예) Clear Buddy
    - p->_mapcount = PAGE_BUDDY_MAPCOUNT_VALUE(-1)
    - new 커널: p->page_type |= 0x80
  - 참고: mm: split page_type out from _mapcount (https://github.com/torvalds/linux/commit/6e292b9be7f4358985ce33ae1f59ab30a8c09e08 )

include/linux/page-flags.h

```
01  /*
02   * PageBuddy() indicates that the page is free and in the buddy system
03   * (see mm/page_alloc.c).
04   */
05  PAGE_TYPE_OPS(Buddy, buddy)
06
07  /*
08   * PageBalloon() is true for pages that are on the balloon page list
09   * (see mm/balloon_compaction.c).
10   */
11  PAGE_TYPE_OPS(Balloon, balloon)
12
13  /*
14   * If kmemcg is enabled, the buddy allocator will set PageKmemcg() on
15   * pages allocated with __GFP_ACCOUNT. It gets cleared on page free.
16   */
17  PAGE_TYPE_OPS(Kmemcg, kmemcg)
18
19  /*
20   * Marks pages in use as page tables.
21   */
```

```
22  PAGE_TYPE_OPS(Table, table)
```

```
1  /*
2   * For pages that are never mapped to userspace (and aren't PageSlab),
3   * page_type may be used.  Because it is initialised to -1, we invert the
4   * sense of the bit, so __SetPageFoo *clears* the bit used for PageFoo, and
5   * __ClearPageFoo *sets* the bit used for PageFoo.  We reserve a few high and
6   * low bits so that an underflow or overflow of page_mapcount() won't be
7   * mistaken for a page type value.
8   */
```

```
01  #define PAGE_TYPE_BASE   0xf0000000
02  /* Reserve               0x0000007f to catch underflows of page_mapcount */
03  #define PAGE_MAPCOUNT_RESERVE   -128
04  #define PG_buddy         0x00000080
05  #define PG_balloon       0x00000100
06  #define PG_kmemcg        0x00000200
07  #define PG_table         0x00000400
08
09  #define PageType(page, flag)                                             \
10          ((page->page_type & (PAGE_TYPE_BASE | flag)) == PAGE_TYPE_BASE)
11
12  static inline int page_has_type(struct page *page)
13  {
14          return (int)page->page_type < PAGE_MAPCOUNT_RESERVE;
15  }
16
17  #define PAGE_TYPE_OPS(uname, lname)                                      \
18  static __always_inline int Page##uname(struct page *page)               \
19  {                                                                        \
20          return PageType(page, PG_##lname);                               \
21  }                                                                        \
22  static __always_inline void __SetPage##uname(struct page *page)         \
23  {                                                                        \
24          VM_BUG_ON_PAGE(!PageType(page, 0), page);                        \
25          page->page_type &= ~PG_##lname;                                  \
26  }                                                                        \
27  static __always_inline void __ClearPage##uname(struct page *page)       \
28  {                                                                        \
29          VM_BUG_ON_PAGE(!Page##uname(page), page);                        \
30          page->page_type |= PG_##lname;                                   \
31  }
```

The above macros create inline functions such as PageBuddy(), __SetPageBuddy(), and __ClearPageBuddy().

# About Page Blocks

## set_pageblock_flags_group()

linux/pageblock-flags.h

```
1   #define set_pageblock_flags_group(page, flags, start_bitidx, end_bitidx)
    \
2           set_pfnblock_flags_mask(page, flags, page_to_pfn(page),
    \
3                           end_bitidx,
    \
4                           (1 << (end_bitidx - start_bitidx + 1)) - 1)
```

## set_pfnblock_flags_mask()

mm/page_alloc.c

```
01  /**
02   * set_pfnblock_flags_mask - Set the requested group of flags for a page
    block_nr_pages block of pages
03   * @page: The page within the block of interest
04   * @flags: The flags to set
05   * @pfn: The target page frame number
06   * @end_bitidx: The last bit of interest
07   * @mask: mask of bits that the caller is interested in
08   */
09  void set_pfnblock_flags_mask(struct page *page, unsigned long flags,
10                                       unsigned long pfn,
11                                       unsigned long end_bitidx,
12                                       unsigned long mask)
13  {
14          struct zone *zone;
15          unsigned long *bitmap;
16          unsigned long bitidx, word_bitidx;
17          unsigned long old_word, word;
18
19          BUILD_BUG_ON(NR_PAGEBLOCK_BITS != 4);
20
21          zone = page_zone(page);
22          bitmap = get_pageblock_bitmap(zone, pfn);
23          bitidx = pfn_to_bitidx(zone, pfn);
24          word_bitidx = bitidx / BITS_PER_LONG;
25          bitidx &= (BITS_PER_LONG-1);
26
27          VM_BUG_ON_PAGE(!zone_spans_pfn(zone, pfn), page);
28
29          bitidx += end_bitidx;
30          mask <<= (BITS_PER_LONG - bitidx - 1);
31          flags <<= (BITS_PER_LONG - bitidx - 1);
32
33          word = ACCESS_ONCE(bitmap[word_bitidx]);
34          for (;;) {
35                  old_word = cmpxchg(&bitmap[word_bitidx], word, (word & ~
    mask) | flags);
36                  if (word == old_word)
37                          break;
38                  word = old_word;
39          }
```

```
40 }
```

## get_pfnblock_flags_mask()

mm/page_alloc.c

```
01  /**
02   * get_pfnblock_flags_mask - Return the requested group of flags for the
     pageblock_nr_pages block of pages
03   * @page: The page within the block of interest
04   * @pfn: The target page frame number
05   * @end_bitidx: The last bit of interest to retrieve
06   * @mask: mask of bits that the caller is interested in
07   *
08   * Return: pageblock_bits flags
09   */
10  unsigned long get_pfnblock_flags_mask(struct page *page, unsigned long p
    fn,
11                                        unsigned long end_bitidx,
12                                        unsigned long mask)
13  {
14          struct zone *zone;
15          unsigned long *bitmap;
16          unsigned long bitidx, word_bitidx;
17          unsigned long word;
18
19          zone = page_zone(page);
20          bitmap = get_pageblock_bitmap(zone, pfn);
21          bitidx = pfn_to_bitidx(zone, pfn);
22          word_bitidx = bitidx / BITS_PER_LONG;
23          bitidx &= (BITS_PER_LONG-1);
24
25          word = bitmap[word_bitidx];
26          bitidx += end_bitidx;
27          return (word >> (BITS_PER_LONG - bitidx - 1)) & mask;
28  }
```

## get_pageblock_bitmap()

mm/page_alloc.c

```
01  /* Return a pointer to the bitmap storing bits affecting a block of page
    s */
02  static inline unsigned long *get_pageblock_bitmap(struct zone *zone,
03                                                     unsigned long pf
    n)
04  {
05  #ifdef CONFIG_SPARSEMEM
06          return __pfn_to_section(pfn)->pageblock_flags;
07  #else
08          return zone->pageblock_flags;
09  #endif /* CONFIG_SPARSEMEM */
10  }
```

Returns a page block bitmap containing the @pfn. (usemap)

- The usemap stores the mobility flags, which are represented in 4 bits.

## pfn_to_bitidx()

mm/page_alloc.c

```
01  static inline int pfn_to_bitidx(struct zone *zone, unsigned long pfn)
02  {
03  #ifdef CONFIG_SPARSEMEM
04          pfn &= (PAGES_PER_SECTION-1);
05          return (pfn >> pageblock_order) * NR_PAGEBLOCK_BITS;
06  #else
07          pfn = pfn - round_down(zone->zone_start_pfn, pageblock_nr_page
    s);
08          return (pfn >> pageblock_order) * NR_PAGEBLOCK_BITS;
09  #endif /* CONFIG_SPARSEMEM */
10  }
```

Returns the bit index from the pageblock for the pfn.

## SECTION_BLOCKFLAGS_BITS

include/linux/mmzone.h

```
1  #define SECTION_BLOCKFLAGS_BITS \
2          ((1UL << (PFN_SECTION_SHIFT - pageblock_order)) * NR_PAGEBLOCK_B
   ITS)
```

Number of pageblock bits per section

- NR_PAGEBLOCK_BITS
  - Number of bits required for pageblock=4
- PFN_SECTION_SHIFT
  - Number of bits required for section length representation – minus number of bits required for page length representation
    - arm64: section length=27 (128M representation) bits – 12 (4KB representation) bits = 15
- e.g. if arm64 has section size=128M, pageblock_order=9
  - SECTION_BLOCKFLAGS_BITS  = 2^(15-9) * 4 bits = 256 bits

# consultation

- Memory Model -1- (Basic) (http://jake.dothome.co.kr/mm-1) | 문c
- Memory Model -2- (mem_map) (http://jake.dothome.co.kr/mem_map) | 문c
- Memory Model -3- (Sparse Memory) (http://jake.dothome.co.kr/sparsemem/) | 문c
- Memory Model -4- (APIs) (http://jake.dothome.co.kr/mem_map_page) | Sentence C – Current post
- ZONE Type (http://jake.dothome.co.kr/zone-types) | Qc
- bootmem_init (http://jake.dothome.co.kr/bootmem_init-64) | Qc
- zone_sizes_init() (http://jake.dothome.co.kr/free_area_init_node/) | Qc
- NUMA -1- (ARM64 initialization) (http://jake.dothome.co.kr/numa-1) | Qc
- build_all_zonelists() (http://jake.dothome.co.kr/build_all_zonelists) | Qc

- An introduction to compound pages (https://lwn.net/Articles/619514/) | LWN.net
- agemap, from the userspace perspective (http://www.kernel.org/doc/Documentation/vm/pagemap.txt) | kernel.org
- Documentation/vm/hugetlbpage.txt (http://www.kernel.org/doc/Documentation/vm/hugetlbpage.txt) | kernel.org

---

**LEAVE A COMMENT**

Your email will not be published. Required fields are marked with *

Comments

name *

email *

Website

WRITE A COMMENT

❮ Memory Model -3- (Sparse Memory) (http://jake.dothome.co.kr/sparsemem/)

__flush_dcache_page() ❯ (http://jake.dothome.co.kr/__flush_dcache_page/)

Munc Blog (2015 ~ 2024)