# Bootmem with bitmap

📅 2016-02-29 (http://jake.dothome.co.kr/bootmem/)    👤 Moon Young-il
(http://jake.dothome.co.kr/author/admin/)    📂 Linux Kernel (http://jake.dothome.co.kr/category/linux/)
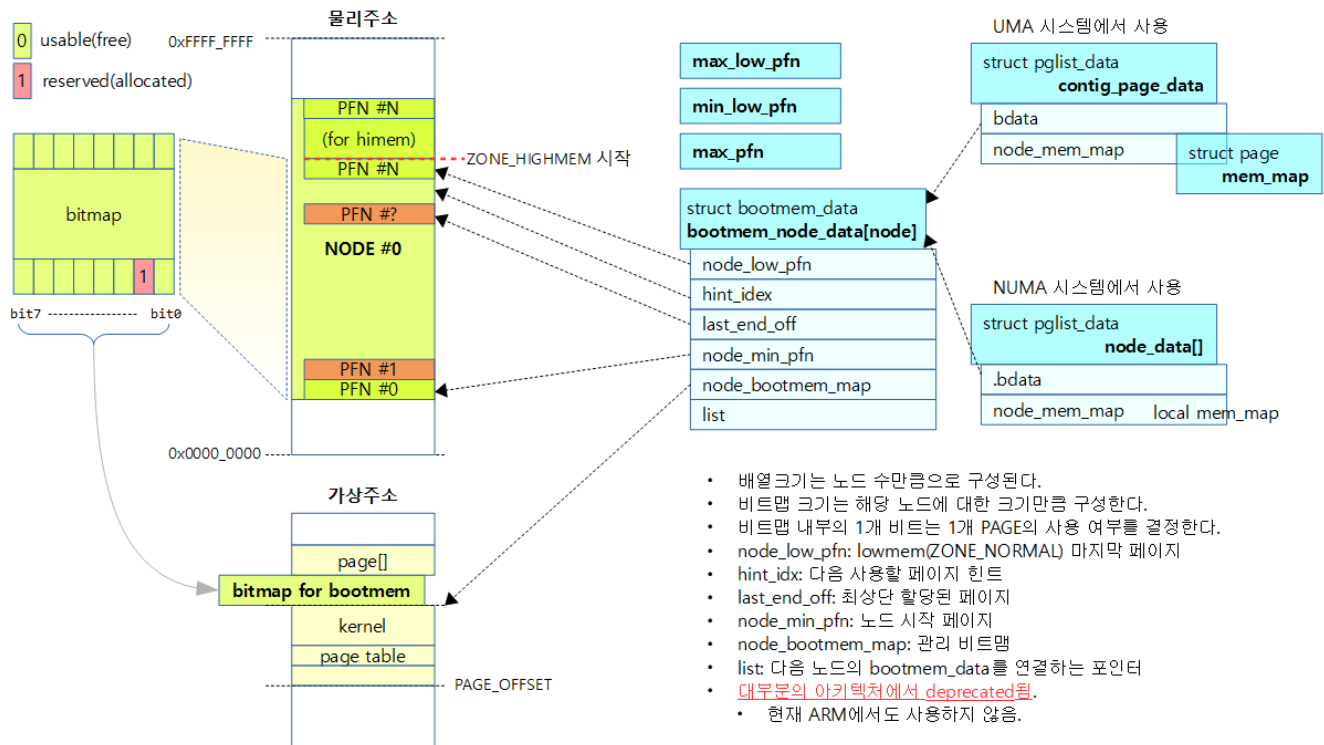
## feature

- A boot-time physical memory allocator and configurator
- Bootmem is simple. During the early part of the kernel bootup process, a low-level memory allocator is used by the kernel.
- Bootmem is responsible for allocating the memory required by the system from the initial MMU (paging) function until the buddy allocator is activated during the early boot process, and is converted to buddy allocation after use.
- Bootmem uses a bitmap to indicate whether or not memory is being used.
- With each architecture, Bootmem is increasingly shifting to a method that uses only memblocks to operate instead of being used in the main kernel.
- If you look at the evolution of memory allocation on x86 systems,
  - 1) very very early allocator (early brk model) [x86] – BIOS "e820"을 사용
  - 2) very early allocator (early_res) -> (memblock) [some generic]
    - Kernel 2.6.35 replaces early_res with memblock
  - 3) early allocator (bootmem) [generic]
  - 4) full buddy allocator
- In ARM kernel v3.14-rc1, arm_bootmem_init() has been removed and an CONFIG_NO_BOOTMEM option has been added.



(http://jake.dothome.co.kr/wp-content/uploads/2016/02/bootmem-1.png)

The figure below shows how bootmem is represented as a bitmap.

(http://jake.dothome.co.kr/wp-content/uploads/2016/02/bootmem-2.png)

# Structure

## struct bootmem_data

include/linux/bootmem.h

```
01  #ifndef CONFIG_NO_BOOTMEM
02  /*
03   * node_bootmem_map is a map pointer - the bits represent all physical
04   * memory pages (including holes) on the node.
05   */
06  typedef struct bootmem_data {
07          unsigned long node_min_pfn;
08          unsigned long node_low_pfn;
09          void *node_bootmem_map;
10          unsigned long last_end_off;
11          unsigned long hint_idx;
12          struct list_head list;
13  } bootmem_data_t;
14
15  extern bootmem_data_t bootmem_node_data[];
16  #endif
```

## struct pglist_data

include/linux/mmzone.h

```
01  /*
02   * The pg_data_t structure is used in machines with CONFIG_DISCONTIGMEM
03   * (mostly NUMA machines?) to denote a higher-level memory zone than the
04   * zone denotes.
05   *
06   * On NUMA machines, each NUMA node would have a pg_data_t to describe
07   * it's memory layout.
08   *
09   * Memory statistics and page replacement data structures are maintained
    on a
```

```
10    * per-zone basis.
11    */
12   struct bootmem_data;
13   typedef struct pglist_data {
14           struct zone node_zones[MAX_NR_ZONES];
15           struct zonelist node_zonelists[MAX_ZONELISTS];
16           int nr_zones;
17   #ifdef CONFIG_FLAT_NODE_MEM_MAP /* means !SPARSEMEM */
18           struct page *node_mem_map;
19   #ifdef CONFIG_PAGE_EXTENSION
20           struct page_ext *node_page_ext;
21   #endif
22   #endif
23   #ifndef CONFIG_NO_BOOTMEM
24           struct bootmem_data *bdata;
25   #endif
26   #ifdef CONFIG_MEMORY_HOTPLUG
27           /*
28            * Must be held any time you expect node_start_pfn, node_present
     _pages
29            * or node_spanned_pages stay constant.  Holding this will also
30            * guarantee that any pfn_valid() stays that way.
31            *
32            * pgdat_resize_lock() and pgdat_resize_unlock() are provided to
33            * manipulate node_size_lock without checking for CONFIG_MEMORY_
     HOTPLUG.
34            *
35            * Nests above zone->lock and zone->span_seqlock
36            */
37           spinlock_t node_size_lock;
38   #endif
39           unsigned long node_start_pfn;
40           unsigned long node_present_pages; /* total number of physical pa
     ges */
41           unsigned long node_spanned_pages; /* total size of physical page
42                                                range, including holes */
43           int node_id;
44           wait_queue_head_t kswapd_wait;
45           wait_queue_head_t pfmemalloc_wait;
46           struct task_struct *kswapd;       /* Protected by
47                                                mem_hotplug_begin/end() */
48           int kswapd_max_order;
49           enum zone_type classzone_idx;
50   #ifdef CONFIG_NUMA_BALANCING
51           /* Lock serializing the migrate rate limiting window */
52           spinlock_t numabalancing_migrate_lock;
53
54           /* Rate limiting time interval */
55           unsigned long numabalancing_migrate_next_window;
56
57           /* Number of pages migrated during the rate limiting time interv
     al */
58           unsigned long numabalancing_migrate_nr_pages;
59   #endif
60   } pg_data_t;
```

## contig_page_data etc.

mm/bootmem.c

```
01   #ifndef CONFIG_NEED_MULTIPLE_NODES
02   struct pglist_data __refdata contig_page_data = {
03           .bdata = &bootmem_node_data[0]
04   };
05   EXPORT_SYMBOL(contig_page_data);
06   #endif
```

```
07
08   unsigned long max_low_pfn;
09   unsigned long min_low_pfn;
10   unsigned long max_pfn;
11
12   bootmem_data_t bootmem_node_data[MAX_NUMNODES] __initdata;
13
14   static struct list_head bdata_list __initdata = LIST_HEAD_INIT(bdata_list);
```

## consultation

- The NO_BOOTMEM patches (http://lwn.net/Articles/382559/) | LWN.net
- mm: Use memblock interface instead of bootmem (https://lwn.net/Articles/576281/) | LWN.net
- Understanding The Linux Virtual Memory Manager | Mel Gorman – 다운로드 (https://www.kernel.org/doc/gorman/pdf/understand.pdf)
- [Linux] bootmem memory allocator (http://egloos.zum.com/studyfoss/v/5026072) | F/OSS

---

**LEAVE A COMMENT**

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

Comments

name *

email *

Website

WRITE A COMMENT

❮ Inline-Assembly (http://jake.dothome.co.kr/inline-assembly/)    Kmap(Pkmap) ❯ (http://jake.dothome.co.kr/kmap/)

Munc Blog (2015 ~ 2023)