

# Zoned Allocator -6- (Watermark)

📅 2016-06-29 (<http://jake.dothome.co.kr/zoned-allocator-watermark/>) 👤 Moon Young-il  
(<http://jake.dothome.co.kr/author/admin/>) 📁 Linux Kernel (<http://jake.dothome.co.kr/category/linux/>)

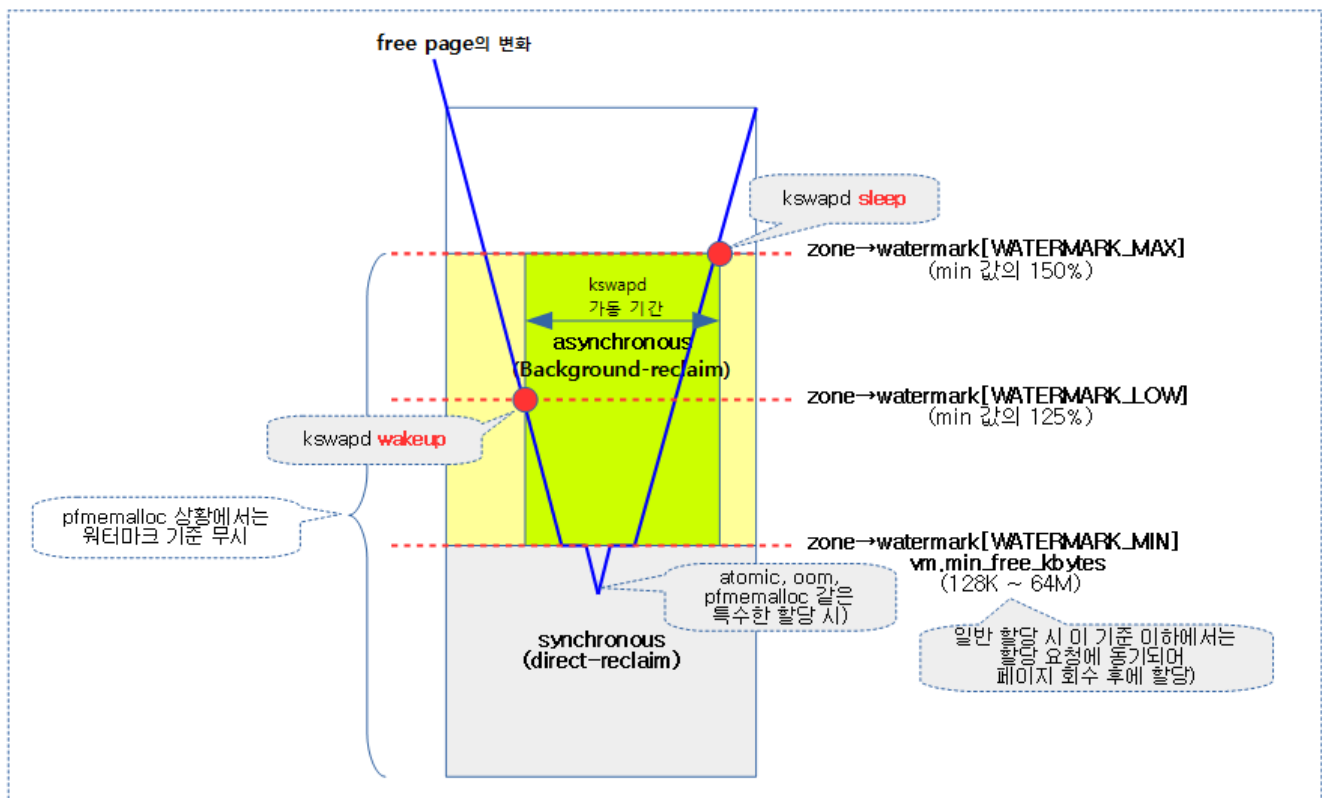
<kernel v5.0>

## Zone Watermark

When attempting to allocate a page in a zone, the remaining free pages are compared with the watermark criteria for each zone to determine whether memory is allocated. These watermarks are used by specifying three watermark reference values as follows, and each of them has the following meanings.

- WATERMARK\_MIN
  - It is equal to the lowest reference point of low memory.
  - `vm.min_free_kbytes` has the number of pages corresponding to the value.
  - If the calculated free page falls below this value, it is not possible to retrieve the page. In this case, if you allow a synchronous direct page reclamation operation, the allocation continues after the page reclaim.
  - Under special circumstances (with the exception of atomic, oom, `pfmemalloc`), page allocation is possible even below this standard.
- WATERMARK\_LOW
  - By default, the number of pages has a value of 125% of the min value. This is the point at which `kswapd` automatically wakes up.
    - Users can change their `watermark_scale_factor` to raise this bar even further.
- WATERMARK\_MAX
  - This is the number of pages with a value of 150% of the min value by default, and this is the point at which `kswapd` is automatically slept.
  - Users can change their `watermark_scale_factor` to raise this bar even further.

The following figure shows the relationship between allocation and page recall based on zone watermark settings.



(<http://jake.dothome.co.kr/wp-content/uploads/2016/06/zone-watermark-1c.png>)

You can check the nr\_free\_pages and min, low, and high watermark values as follows.

```
$ cat /proc/zoneinfo
Node 0, zone DMA32
per-node stats
(...생략...)
pages free      633680
    min         5632      <-----
    low         7040      <-----
    high        8448      <-----
spanned 786432
present 786432
managed 765771
protection: (0, 0, 0)
nr_free_pages 633680      <-----
(...생략...)
```

## managed\_pages

The number of pages that the buddy system manages for each zone and the initial values in the DMA and normal zones are as follows:

- managed\_pages = present\_pages – the number of pages that were reserved just before the buddy system started
- managed\_pages = spanned\_pages – absent\_pages – (memmap\_pages + dma\_reserve + ...)
  - present\_pages = spanned\_pages – absent\_pages

You can check the managed\_pages values for each zone as follows.

```
$ cat /proc/zoneinfo
Node 0, zone      DMA32
  per-node stats
  (...생략...)
    pages free      633680
      min          5632
      low          7040
      high          8448
    spanned 786432
    present 786432
    managed 765771      <-----
    protection: (0, 0, 0)
    nr_free_pages 633680
  (...생략...)
```

## Quick Watermark Comparison

### zone\_watermark\_fast()

mm/page\_alloc.c

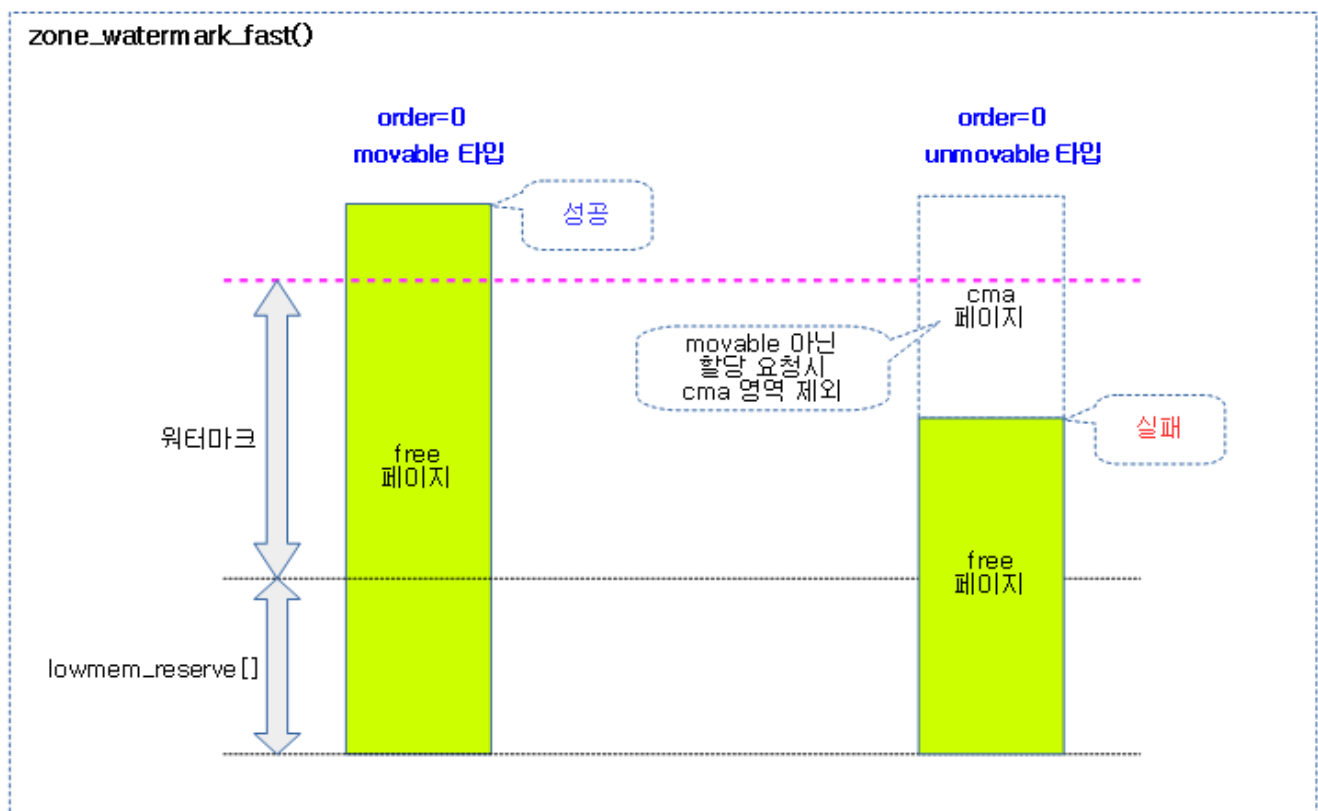
```
01 static inline bool zone_watermark_fast(struct zone *z, unsigned int orde
02 r,
03     unsigned long mark, int classzone_idx, unsigned int allo
04 c_flags)
05 {
06     long free_pages = zone_page_state(z, NR_FREE_PAGES);
07     long cma_pages = 0;
08     #ifdef CONFIG_CMA
09     /* If allocation can't use CMA areas don't use free CMA pages */
10     if (!(alloc_flags & ALLOC_CMA))
11         cma_pages = zone_page_state(z, NR_FREE_CMA_PAGES);
12     #endif
13     /*
14      * Fast check for order-0 only. If this fails then the reserves
15      * need to be calculated. There is a corner case where the check
16      * passes but only the high-order atomic reserve are free. If
17      * the caller is !atomic then it'll uselessly search the free
18      * list. That corner case is then slower but it is harmless.
19      */
20     if (!order && (free_pages - cma_pages) > mark + z->lowmem_reserv
21 e[classzone_idx])
22         return true;
23     return __zone_watermark_ok(z, order, mark, classzone_idx, alloc_
24 flags,
25                             free_pages);
```

It quickly determines whether the requested @z zone can secure @order pages from a free page that is higher than the @mark standard. Judgment logic has been added exclusively for order 0 allocation requests. (true=OK above the watermark standard, false=NOT-OK below the watermark standard)

- In line 4 of code, we find the free page for that zone.

- If the page request is not movable in line 9~10 of the code, the pages in the cma area will not be available. First, we know how many pages are using the CMA area.
- In line 19~20 of the code, if the request is to assign order 0, you can make a quick judgment. Returns true if the free page, except for the CMA page, exceeds the watermark threshold plus the lowmem reserve area.
- In line 22~23 of the code, order returns whether the allocation request of 1 or higher exceeds the watermark criterion.

The following figure shows how the order 0 page allocation request is judged by the watermark criterion for quick assignment.



([http://jake.dothome.co.kr/wp-content/uploads/2016/06/zone\\_watermark\\_fast-1a.png](http://jake.dothome.co.kr/wp-content/uploads/2016/06/zone_watermark_fast-1a.png))

## Comparison of Default Watermarks

### zone\_watermark\_ok()

mm/page\_alloc.c

```

1 | bool zone_watermark_ok(struct zone *z, unsigned int order, unsigned long
  | mark,
2 |                               int classzone_idx, int alloc_flags)
3 | {
4 |     return __zone_watermark_ok(z, order, mark, classzone_idx, alloc_
  | flags,
5 |                               zone_page_state(z, NR_FREE_PAGE
6 | S));
  | }

```

In order to handle the order page in that zone, make sure that the free page is sufficient for the watermark value.

## \_\_zone\_watermark\_ok()

mm/page\_alloc.c

```

1  /*
2   * Return true if free base pages are above 'mark'. For high-order check
   s it
3   * will return true of the order-0 watermark is reached and there is at
   least
4   * one free page of a suitable size. Checking now avoids taking the zone
   lock
5   * to check in the allocation paths if no pages are free.
6   */

01 bool __zone_watermark_ok(struct zone *z, unsigned int order, unsigned lo
ng mark,
02                          int classzone_idx, unsigned int alloc_flags,
03                          long free_pages)
04 {
05     long min = mark;
06     int o;
07     const bool alloc_harder = (alloc_flags & (ALLOC_HARDER|ALLOC_00
M));
08
09     /* free_pages may go negative - that's OK */
10     free_pages -= (1 << order) - 1;
11
12     if (alloc_flags & ALLOC_HIGH)
13         min -= min / 2;
14
15     /*
16      * If the caller does not have rights to ALLOC_HARDER then subtr
17 act
18 the
19     * the high-atomic reserves. This will over-estimate the size of
20     * atomic reserve but it avoids a search.
21     */
22     if (likely(!alloc_harder)) {
23         free_pages -= z->nr_reserved_highatomic;
24     } else {
25         /*
26          * OOM victims can try even harder than normal ALLOC_HAR
27 DER
28 in
29 it
30 ved.
31         * makes during the free path will be small and short-li
32 ved.
33         */
34         if (alloc_flags & ALLOC_OOM)
35             min -= min / 2;
36         else
37             min -= min / 4;
38     }
39
40 #ifdef CONFIG_CMA
41     /* If allocation can't use CMA areas don't use free CMA pages */
42     if (!(alloc_flags & ALLOC_CMA))
43         free_pages -= zone_page_state(z, NR_FREE_CMA_PAGES);
44 #endif

```

```

41
42     /*
43     * Check watermarks for an order-0 allocation request. If these
44     * are not met, then a high-order request also cannot go ahead
45     * even if a suitable page happened to be free.
46     */
47     if (free_pages <= min + z->lowmem_reserve[classzone_idx])
48         return false;
49
50     /* If this is an order-0 request then the watermark is fine */
51     if (!order)
52         return true;
53
54     /* For a high-order request, check at least one suitable page is
55     free */
56     for (o = order; o < MAX_ORDER; o++) {
57         struct free_area *area = &z->free_area[o];
58         int mt;
59
60         if (!area->nr_free)
61             continue;
62
63         for (mt = 0; mt < MIGRATE_PCPTYPES; mt++) {
64             if (!list_empty(&area->free_list[mt]))
65                 return true;
66         }
67     #ifdef CONFIG_CMA
68         if ((alloc_flags & ALLOC_CMA) &&
69             !list_empty(&area->free_list[MIGRATE_CMA])) {
70             return true;
71         }
72     #endif
73     if (alloc_harder &&
74         !list_empty(&area->free_list[MIGRATE_HIGHATOMIC]))
75         return true;
76     }
77     return false;
78 }

```

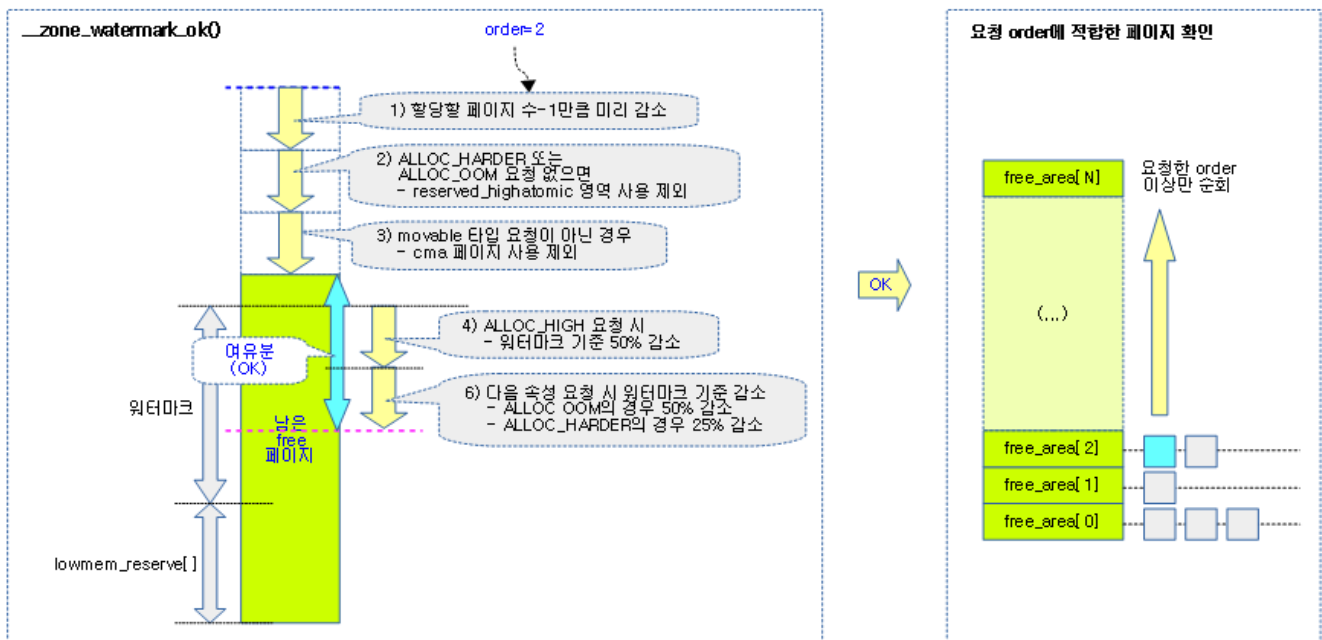
In order to handle the order page in that zone, make sure that the free page is sufficient for the watermark value.

- In line 5 of the code, assign the watermark value to the min value.
- In line 7 of the code, it is decided whether the watermark standard should be lowered as much as possible to determine whether the allocation should be made in the alloc\_harder.
  - Allocation in gfp\_atomic requests and OOM situations requires the system to complete the allocation as much as possible.
- In line 10 of code, the number of pages requested to allocate free pages is reduced by -1.
- In line 12~13 of the code, the ALLOC\_HIGH assignment flag is set in the case of an atomic request, and the watermark criteria are halved to allow for maximum allocation.
- Unless it is alloc\_harder in line 20~21 of the code, the reserved\_highatomic pages are omitted to prevent use.
- In line 22~33 of the code, the watermark is lowered further in the case of alloc\_harder situation, which is 25% lower in the case of an atomic request and 50% in the case of an OOM situation.
- If the page allocation request is not movable on lines 38~39 of the code, the cma area cannot be used. Therefore, the pages in the CMA area are excluded from the free page count first.

- In line 47~48 of the code, if the free page is below the watermark threshold and is low on memory, including the lowmem reserve area, it will return false.
- On lines 51~52 of the code, it returns true for an order 0 assignment request.
- In line 55~76 of the code, for requests with order 1 or more, the actual free list is searched and returned after checking whether it is assignable. [request order, MAX\_ORDER) traverses the free\_area in the range, and returns true because it is assignable if it is included in the following 3 criteria:
  - If there is a free page in the list of type unmovable, movable, or reclaimable.
  - If the request is a movable page that can use the CMA area, if there is a free page in the list of type CMA.
  - In alloc\_harder situation, if there is a free page in a list of type highatomic

The following figure shows checking whether a free page can be assigned according to the zone watermark setting.

- After checking whether there is enough space for the remaining free pages as shown in the figure on the left,
- Search for free\_area[], as shown in the figure on the right, to see if it is finally assignable.
- The zone index used in lowmem\_reserve[] is the preferred zone index that was initially requested.



([http://jake.dothome.co.kr/wp-content/uploads/2016/06/zone\\_watermark\\_ok-1.png](http://jake.dothome.co.kr/wp-content/uploads/2016/06/zone_watermark_ok-1.png))

## Accurate Watermark Comparison

### zone\_watermark\_ok\_safe()

page\_alloc.c

```
01 | bool zone_watermark_ok_safe(struct zone *z, unsigned int order,
02 |                             unsigned long mark, int classzone_idx)
03 | {
04 |     long free_pages = zone_page_state(z, NR_FREE_PAGES);
```

```

05
06     if (z->percpu_drift_mark && free_pages < z->percpu_drift_mark)
07         free_pages = zone_page_state_snapshot(z, NR_FREE_PAGES);
08
09     return __zone_watermark_ok(z, order, mark, classzone_idx, 0,
10                                free_pag
11     es);
12 }

```

In order to process the order page in the zone, make sure that the correctly calculated free page is sufficient for the watermark value.

- In line 4 of the code, we first get the approximate number of free pages.
- If the number of free pages from code lines 6~7 is less than `percpu_drift_mark`, the free pages are calculated correctly.
  - `percpu_drift_mark` value is set in the `refresh_zone_stat_thresholds()` function.
- In line 9~10 of the code, compare the number of free pages with the watermark and return whether it can be assigned or not.

## Zone Counter (STAT)

The zone counters (stats) have two separate zone counters and two per-CPU counters for performance. If only an approximate value of the zone is required, the zone counter value is returned as it is for performance, and the zone counter and per-cpu counter are added only if the correct calculation is required.

### **percpu\_drift\_mark**

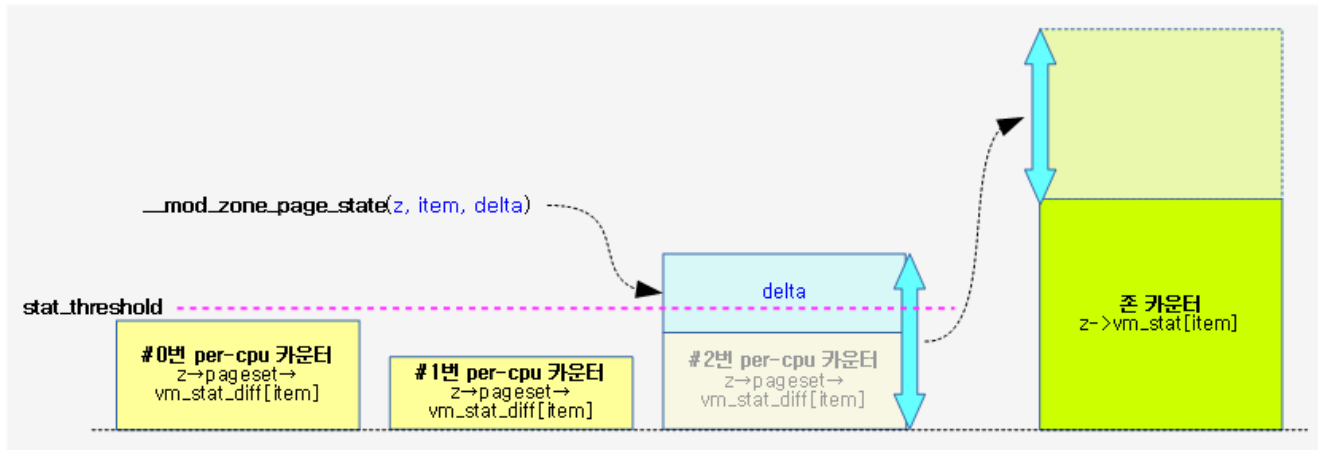
In kernel memory management, a routine is often used to determine memory shortage by reading the number of free pages remaining and comparing it with the watermark. However, in order to read the correct free page value, both the zone counter and the per-cpu counter must be read and added, and this calculation decreases the performance, so the performance is maintained by reading only the zone counter value and comparing it with the `percpu_drift_mark` value set to a certain size larger than the high watermark and inducing more accurate calculations only when it is below this value.

### **stat\_threshold**

When the counter increases or decreases, the per-CPU counter is increased or decreased first for performance. For the accuracy of the zone counter, the value of the per-cpu counter is transferred to the zone counter and added only when the value exceeds a certain standard `stat_threshold`.

The following figure shows how when the zone counter of the CPU changes, it is moved to the zone counter if it exceeds `stat_threshold`.





([http://jake.dothome.co.kr/wp-content/uploads/2016/06/mod\\_zone\\_page\\_state-1a.png](http://jake.dothome.co.kr/wp-content/uploads/2016/06/mod_zone_page_state-1a.png))

You can check the `stat_threshold` values as follows:

```
$ cat /proc/zoneinfo
Node 0, zone    DMA32
per-node stats
(...생략...)
pagesets
cpu: 0
    count: 143
    high: 378
    batch: 63
vm stats threshold: 24    <-----
cpu: 1
    count: 285
    high: 378
    batch: 63
vm stats threshold: 24    <-----
(...생략...)
```

### zone\_page\_state()

include/linux/vmstat.h

```
01 | static inline unsigned long zone_page_state(struct zone *zone,
02 |                                             enum zone_stat_item item)
03 | {
04 |     long x = atomic_long_read(&zone->vm_stat[item]);
05 |     #ifdef CONFIG_SMP
06 |         if (x < 0)
07 |             x = 0;
08 |     #endif
09 |     return x;
10 | }
```

Returns the approximate VM stat value for the zone's @item.

- Returns only the vm stat value of the zone.

### zone\_page\_state\_snapshot()

include/linux/vmstat.h

```

1  /*
2  * More accurate version that also considers the currently pending
3  * deltas. For that we need to loop over all cpus to find the current
4  * deltas. There is no synchronization so the result cannot be
5  * exactly accurate either.
6  */

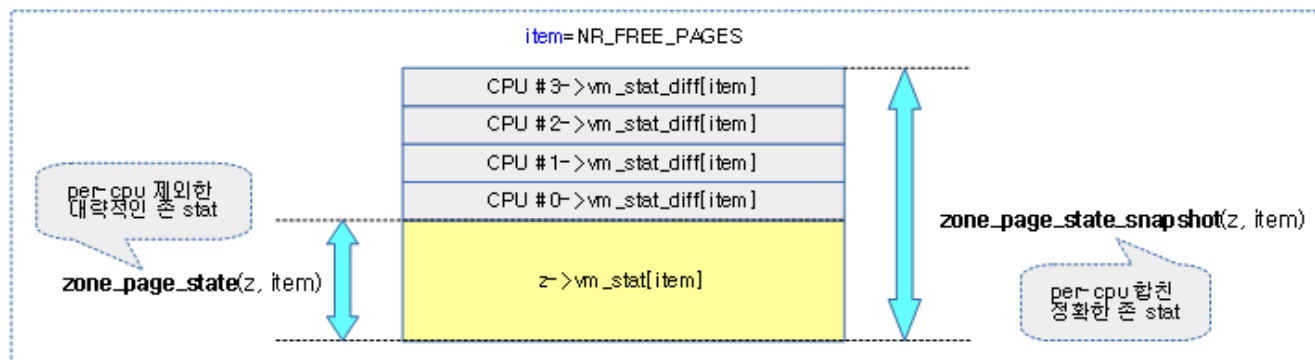
01 static inline unsigned long zone_page_state_snapshot(struct zone *zone,
02                                                       enum zone_stat_item item)
03 {
04     long x = atomic_long_read(&zone->vm_stat[item]);
05
06     #ifdef CONFIG_SMP
07         int cpu;
08         for_each_online_cpu(cpu)
09             x += per_cpu_ptr(zone->pageset, cpu)->vm_stat_diff[ite
10 m];
11
12         if (x < 0)
13             x = 0;
14     #endif
15     return x;

```

Accurately calculates the vm stat value corresponding to the zone's @item.

- Use the following formula to return the exact value:
  - vm stat value in zone + all per-cpu vm stat value

The following figure shows the difference between the zone\_page\_state\_snapshot() function and the zone\_page\_state() function.

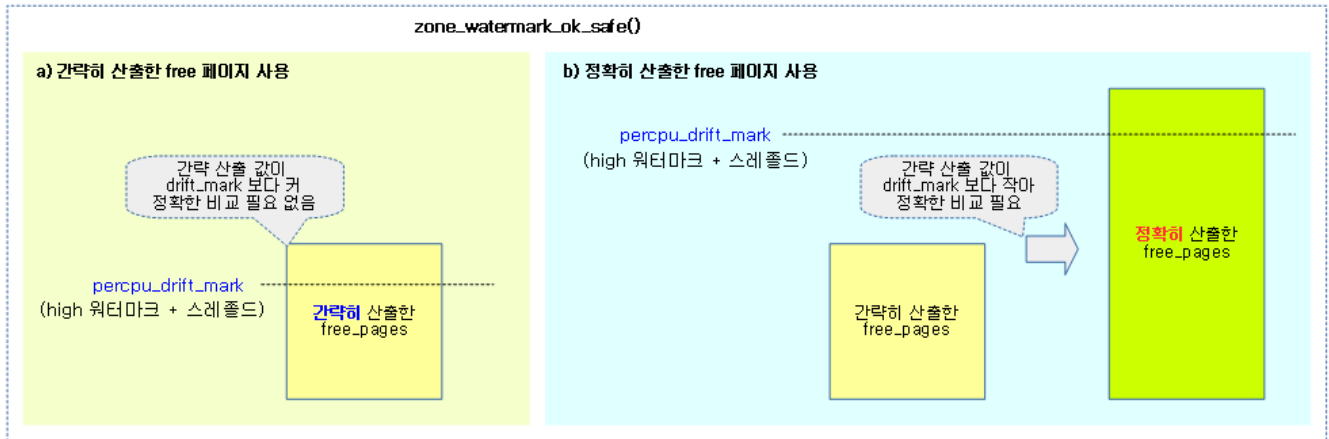


([http://jake.dothome.co.kr/wp-content/uploads/2016/06/zone\\_page\\_state\\_snapshot-1a.png](http://jake.dothome.co.kr/wp-content/uploads/2016/06/zone_page_state_snapshot-1a.png))

### What's the difference between the zone\_watermark\_ok() function and the zone\_watermark\_ok\_safe() function?

- zone\_watermark\_ok() function
  - Use the free page value, which is roughly calculated for the number of free pages to compare with the watermark.
- zone\_watermark\_ok\_safe() function
  - Use the free page value calculated exactly for the number of free pages to be compared with the watermark.

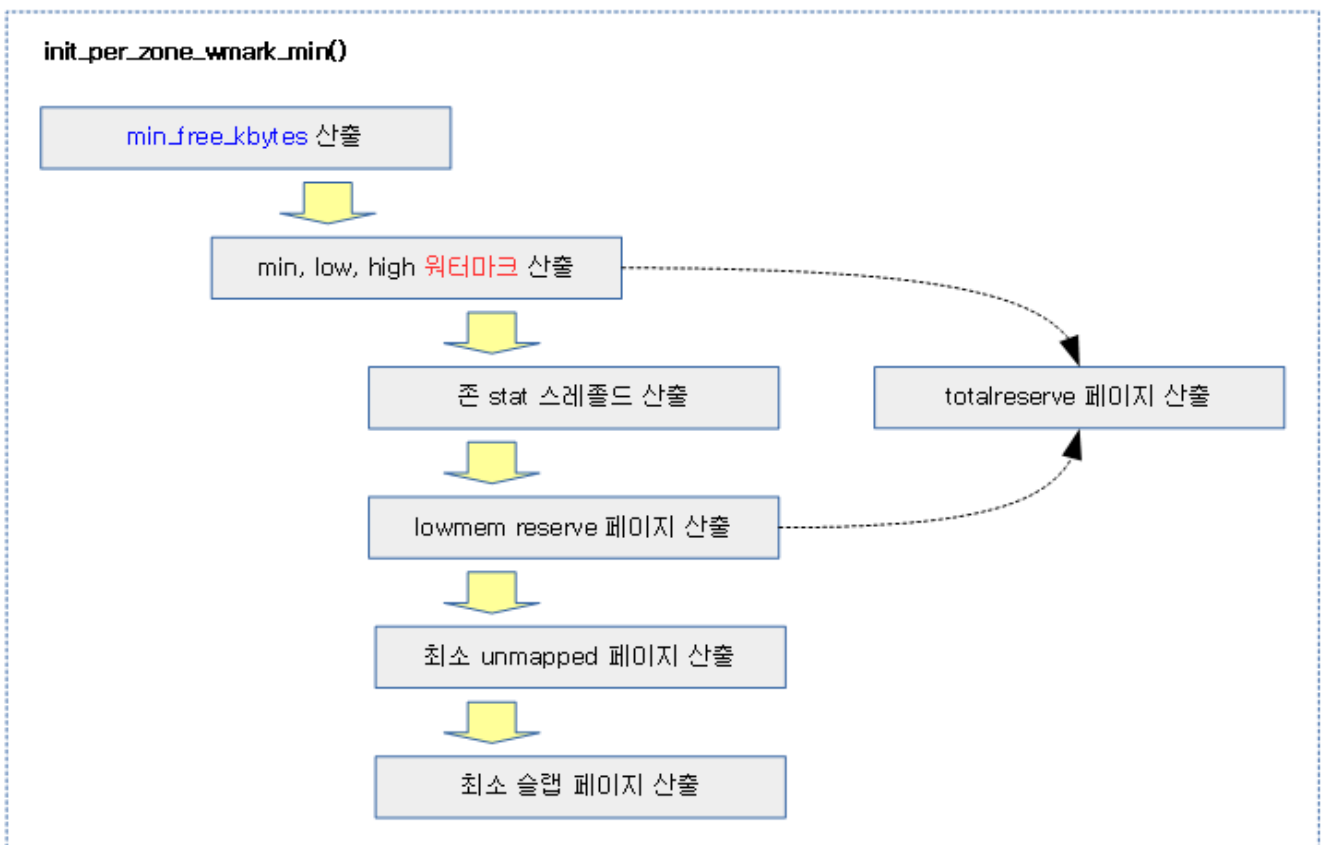
Below is an example of the use of the `zone_watermark_ok_safe()` function using the zone stat thread.



([http://jake.dothome.co.kr/wp-content/uploads/2016/06/zone\\_watermark\\_ok\\_safe-1a.png](http://jake.dothome.co.kr/wp-content/uploads/2016/06/zone_watermark_ok_safe-1a.png))

## Reset watermark and related settings

Resetting the watermark and related settings is done as follows.



([http://jake.dothome.co.kr/wp-content/uploads/2016/06/init\\_per\\_zone\\_wmark\\_min-2a.png](http://jake.dothome.co.kr/wp-content/uploads/2016/06/init_per_zone_wmark_min-2a.png))

## Set the watermark min value

`init_per_zone_wmark_min()`

## mm/page\_alloc.c

```

01  /*
02  *  Initialise min_free_kbytes.
03  *
04  *  For small machines we want it small (128k min).  For large machines
05  *  we want it large (64MB max).  But it is not linear, because network
06  *  bandwidth does not increase linearly with machine size.  We use
07  *
08  *      min_free_kbytes = 4 * sqrt(lowmem_kbytes), for better accuracy:
09  *      min_free_kbytes = sqrt(lowmem_kbytes * 16)
10  *
11  *  which yields
12  *
13  *  16MB:      512k
14  *  32MB:      724k
15  *  64MB:      1024k
16  *  128MB:     1448k
17  *  256MB:     2048k
18  *  512MB:     2896k
19  *  1024MB:    4096k
20  *  2048MB:    5792k
21  *  4096MB:    8192k
22  *  8192MB:    11584k
23  *  16384MB:   16384k
24  */

01  int __meminit init_per_zone_wmark_min(void)
02  {
03      unsigned long lowmem_kbytes;
04      int new_min_free_kbytes;
05
06      lowmem_kbytes = nr_free_buffer_pages() * (PAGE_SIZE >> 10);
07      new_min_free_kbytes = int_sqrt(lowmem_kbytes * 16);
08
09      if (new_min_free_kbytes > user_min_free_kbytes) {
10          min_free_kbytes = new_min_free_kbytes;
11          if (min_free_kbytes < 128)
12              min_free_kbytes = 128;
13          if (min_free_kbytes > 65536)
14              min_free_kbytes = 65536;
15      } else {
16          pr_warn("min_free_kbytes is not updated to %d because us
er defined value %d is pref
ered\n",
17                  new_min_free_kbytes, user_min_free_kbyte
s);
18      }
19      setup_per_zone_wmarks();
20      refresh_zone_stat_thresholds();
21      setup_per_zone_lowmem_reserve();
22
23
24  #ifdef CONFIG_NUMA
25      setup_min_unmapped_ratio();
26      setup_min_slab_ratio();
27  #endif
28
29      return 0;
30  }
31  core_initcall(init_per_zone_wmark_min)

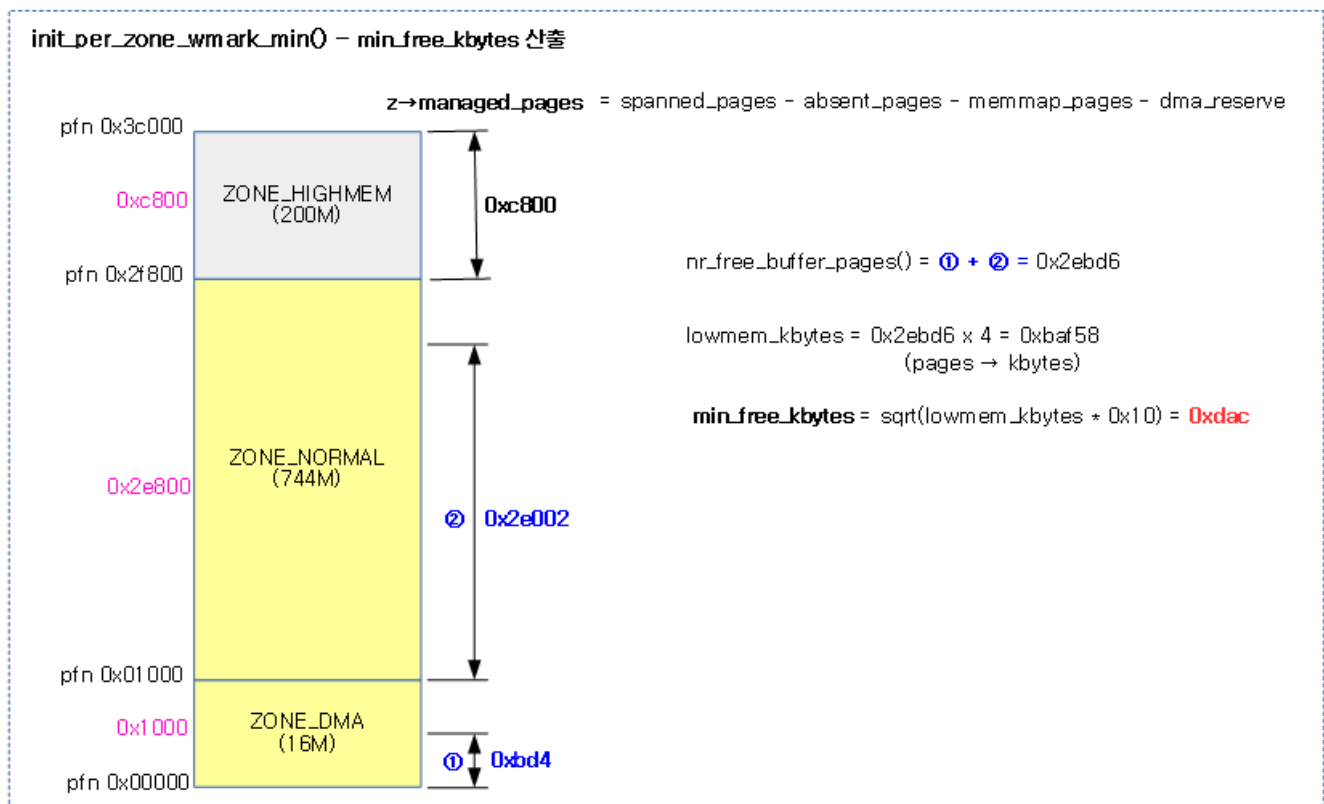
```

It uses lowmem's available page size to calculate the min\_free\_kbytes value used as the basis for the watermark. In addition, it calculates watermarks, lowmem\_reserve, totalreserve, min unpapped, min slap pages, etc. for each zone.

- In line 6~18 of the code, the number of available pages in lowmem is used to calculate the min\_free\_kbytes value used as the watermark reference value.
  - Lowmem Available Page Calculation Formula
    - managed\_pages – High Watermark Page
  - min\_free\_kbytes Calculation Formula
    - $\text{sqrt}((\text{Number of available pages in lowmem} \gg 10) * 16)$ 
      - 예)  $\text{sqrt}(512\text{M} * 0\text{x}10) \rightarrow 2\text{M}$
    - Calculated values are limited to the range of 128 ~ 65536 (range of 128K ~ 64M)
  - In addition to initializing the kernel, this initialization routine is performed whenever the memory is hot-add/del using hotplug memory.
  - When the kernel is initialized for the first time, the initial value of each zone watermark is 0.
- In line 19 of the code, it calculates the watermark value and totalreserve\_pages value for each zone.
- On line 20 of code, we calculate the stat thread and percpu\_drift\_mark for each zone.
- Line 21 of the code yields the zone->lowmem\_reserve and totalreserve pages.
- Yield at least ummapped pages at line 24 of code.
- Yield a minimum of slab pages on line 25 of code.

The following figure shows how the global min\_free\_kbytes value is calculated.

- It is calculated by multiplying the available pages (managed – high watermark) of lowmem (dma, normal) by 16 and root to the square.
  - On the first bootup, the high watermark value is 0.
- In the case of ARM64, the dma32 zone is used, and the highmem zone is not used.



([http://jake.dothome.co.kr/wp-content/uploads/2016/06/init\\_per\\_zone\\_wmark\\_min-1a.png](http://jake.dothome.co.kr/wp-content/uploads/2016/06/init_per_zone_wmark_min-1a.png))

## Calculate the number of lowmem available pages

### nr\_free\_buffer\_pages()

mm/page\_alloc.c

```

1  /**
2   * nr_free_buffer_pages - count number of pages beyond high watermark
3   *
4   * nr_free_buffer_pages() counts the number of pages which are beyond th
5   * e high
6   * watermark within ZONE_DMA and ZONE_NORMAL.
7   */
8
9  unsigned long nr_free_buffer_pages(void)
10 {
11     return nr_free_zone_pages(gfp_zone(GFP_USER));
12 }
13 EXPORT_SYMBOL_GPL(nr_free_buffer_pages);

```

Returns the number of available pages in the current node's zonelist for zones below the normal zone.

- Note that the first time it is called to configure the kernel, the high watermark is 0.
- Number of pages available
  - managed\_pages – High Watermark Page

### nr\_free\_zone\_pages()

mm/page\_alloc.c

```

01  /**
02   * nr_free_zone_pages - count number of pages beyond high watermark
03   * @offset: The zone index of the highest zone
04   *
05   * nr_free_zone_pages() counts the number of counts pages which are beyo
06   * nd the
07   * high watermark within all zones at or below a given zone index. For
08   * each
09   * zone, the number of pages is calculated as:
10   *
11   * nr_free_zone_pages = managed_pages - high_pages
12   */
13
14  static unsigned long nr_free_zone_pages(int offset)
15  {
16     struct zoneref *z;
17     struct zone *zone;
18
19     /* Just pick one node, since fallback list is circular */
20     unsigned long sum = 0;
21
22     struct zonelist *zonelist = node_zonelist(numa_node_id(), GFP_KERNEL);
23
24     for_each_zone_zonelist(zone, z, zonelist, offset) {
25         unsigned long size = zone_managed_pages(zone);
26         unsigned long high = high_wmark_pages(zone);
27         if (size > high)
28             sum += size - high;
29     }
30 }

```

```
return sum;
```

- Number of pages available
  - `managed_pages` – High Watermark Page

The spacing between watermarks uses 25% of the watermark min value by default, and users can change the watermark ratio to increase kswapd efficiency.

/proc/sys/vm/watermark\_scale\_factor and has an initial value of 10. When the initial value is 10, it means that it is the value of 0.1% of the memory, and when the maximum value is 1000, it is the value of 10% of the memory. In this case, the maximum interval between min and low and low and high is possible up to 10% of the memory (managed\_pages).

- In kernel v4.6-rc1, the threshold of the watermark low and high values has been added to increase the efficiency of kswapd on systems with large memory.

Using `min_free_kbytes` values determined by the amount of memory, min, low, and high watermarks are calculated for each zone.

## mm/page\_alloc.c

```
1  /**
2   * setup_per_zone_wmarks - called when min_free_kbytes changes
3   * or when memory is hot-{added|removed}
4   *
5   * Ensures that the watermark[min,low,high] values for each zone are set
6   * correctly with respect to min_free_kbytes.
7   */

1  void setup_per_zone_wmarks(void)
2  {
3      mutex_lock(&zonelists_mutex);
4      __setup_per_zone_wmarks();
5      mutex_unlock(&zonelists_mutex);
6  }
```

[jake.dothome.co.kr/zoned-allocator-watermark/](http://jake.dothome.co.kr/zoned-allocator-watermark/)

**\_\_setup\_per\_zone\_wmarks()**

mm/page\_alloc.c

```

01 static void __setup_per_zone_wmarks(void)
02 {
03     unsigned long pages_min = min_free_kbytes >> (PAGE_SHIFT - 10);
04     unsigned long lowmem_pages = 0;
05     struct zone *zone;
06     unsigned long flags;
07
08     /* Calculate total number of !ZONE_HIGHMEM pages */
09     for_each_zone(zone) {
10         if (!is_highmem(zone))
11             lowmem_pages += zone_managed_pages(zone);
12     }
13
14     for_each_zone(zone) {
15         u64 tmp;
16
17         spin_lock_irqsave(&zone->lock, flags);
18         tmp = (u64)pages_min * zone_managed_pages(zone);
19         do_div(tmp, lowmem_pages);
20         if (is_highmem(zone)) {
21             /*
22              * __GFP_HIGH and PF_MEMALLOC allocations usuall
23              * need highmem pages, so cap pages_min to a sma
24              * value here.
25              *
26              * The WMARK_HIGH-WMARK_LOW and (WMARK_LOW-WMARK
27              * _MIN)
28              * deltas control asynch page reclaim, and so sh
29              * ould
30              * not be capped for highmem.
31              */
32             unsigned long min_pages;
33             min_pages = zone_managed_pages(zone) / 1024;
34             min_pages = clamp(min_pages, SWAP_CLUSTER_MAX, 1
35             28UL);
36             zone->_watermark[WMARK_MIN] = min_pages;
37         } else {
38             /*
39              * If it's a lowmem zone, reserve a number of pa
40              * ges
41              * proportionate to the zone's size.
42              */
43             zone->_watermark[WMARK_MIN] = tmp;
44         }
45
46         /*
47          * Set the kswapd watermarks distance according to the
48          * scale factor in proportion to available memory, but
49          * ensure a minimum size on small systems.
50          */
51         tmp = max_t(u64, tmp >> 2,
52             mult_frac(zone_managed_pages(zone),
53             watermark_scale_factor, 10000));
54         zone->_watermark[WMARK_LOW] = min_wmark_pages(zone) + t
55         mp;
56         zone->_watermark[WMARK_HIGH] = min_wmark_pages(zone) + t
57         mp * 2;
58         zone->watermark_boost = 0;
59
60         spin_unlock_irqrestore(&zone->lock, flags);

```



```

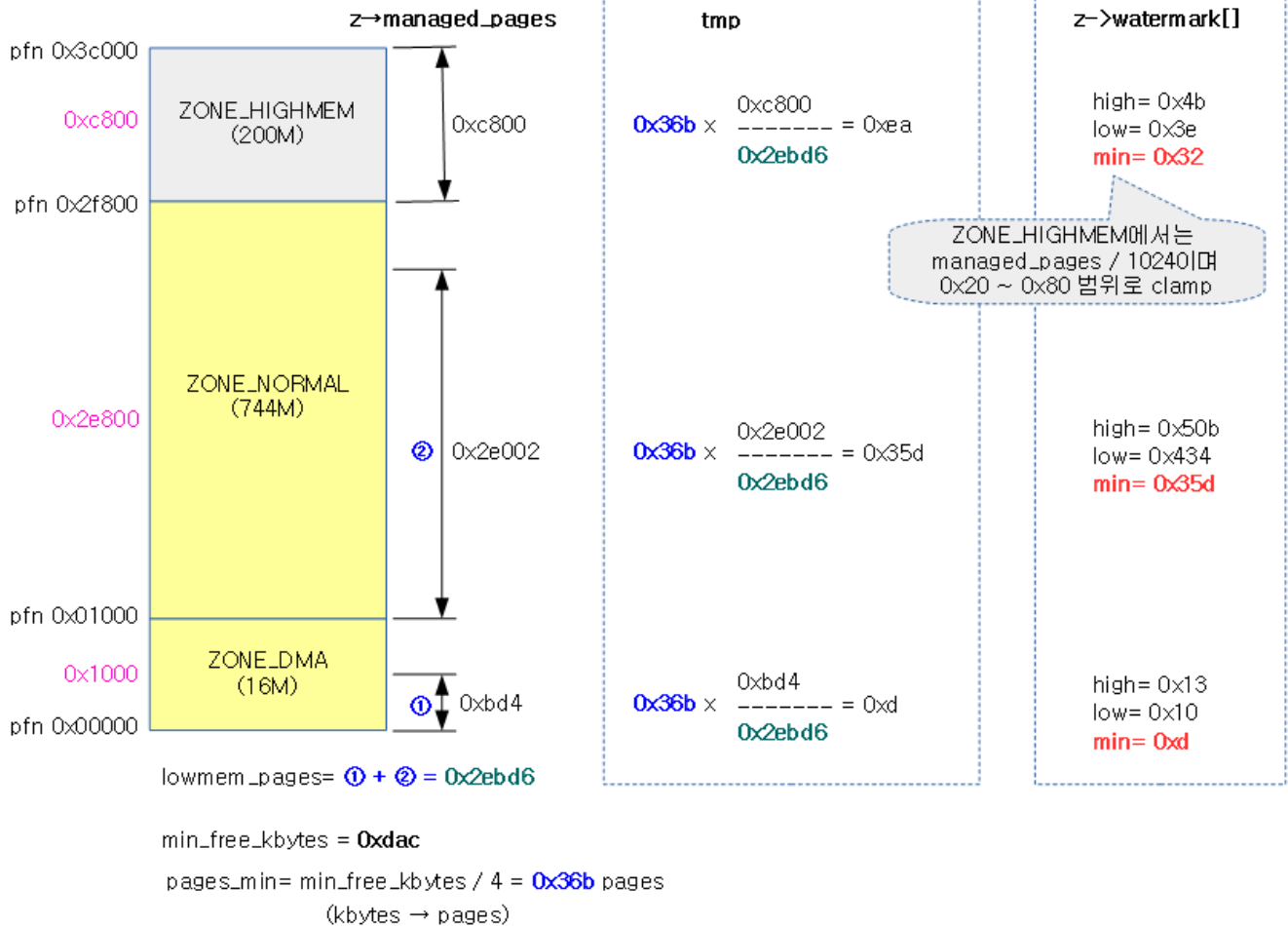
57     }
58
59     /* update totalreserve_pages */
60     calculate_totalreserve_pages();
61 }

```

Calculate low, min, and high watermarks for each zone and update the TotalReserve page.

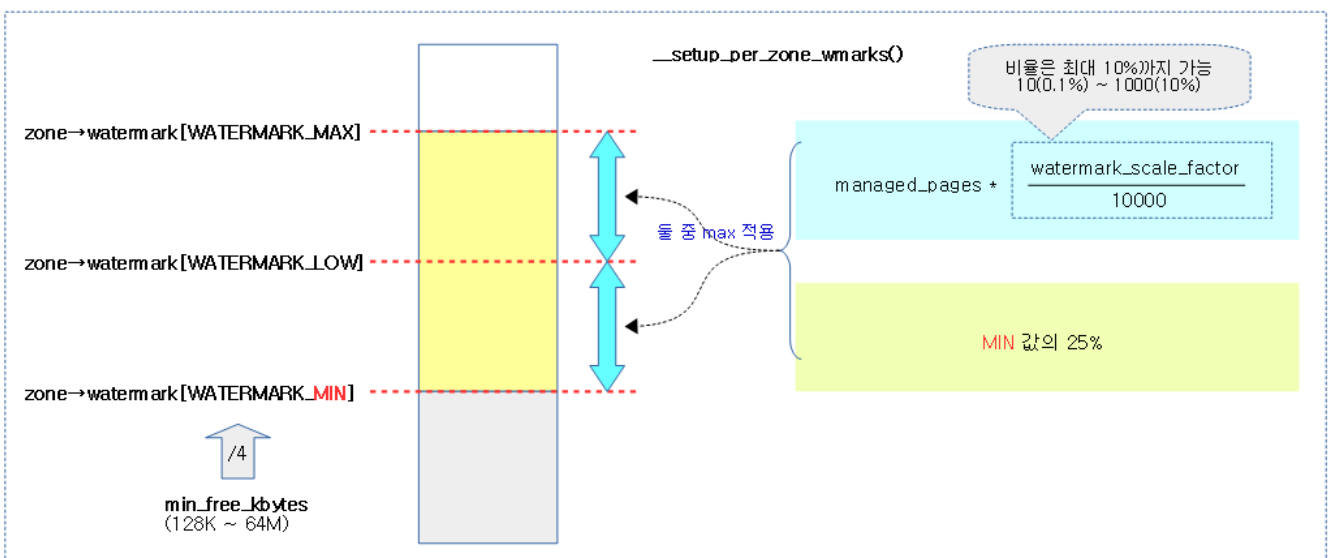
- In line 3 of code, replace the global min\_free\_kbytes value in terms of pages and assign it to the pages\_min.
- In line 9~12 of the code, sum the managed\_pages of the lowmem area to get the lowmem\_pages.
- In code lines 14~19, go through each zone and set the pages\_min (min\_free\_kbytes to pages) to the number of pages equal to the percentage of the current zone in the TMP.
- In line 20~41 of the code, if it is a highmem zone, the actual managed\_pages divided by 1024 is limited to the range of 32~128 and the value is stored in the min watermark, and if it is not a highmem zone, the tmp value calculated above is set in the min watermark.
- In code lines 48~53, the spacing between each watermark is determined by the greater of the following two calculated values.
  - 25% of the initial calculated min value
  - managed\_pages \* Watermark Ratio (0.1 ~ 10%)
- In line 54 of code, reset the Watermark Boost value to 0.
- Since the high watermark value has been updated in line 60 of the code, it is updated by calculating the number of pages in totalreserve.
  - Whenever the high watermark and lowmem reserve pages are recalculated, the totalreserve value is also updated.

The following figure shows the min\_free\_kbytes value and the available pages (managed\_pages – high watermarks) in the lowmem area and the min, low, and high watermarks in the proportions of each zone.



([http://jake.dothome.co.kr/wp-content/uploads/2016/06/setup\\_per\\_zone\\_wmarks-1.png](http://jake.dothome.co.kr/wp-content/uploads/2016/06/setup_per_zone_wmarks-1.png))

The following figure shows the largest of the two calculated values on the right as the watermark spacing value.



([http://jake.dothome.co.kr/wp-content/uploads/2016/06/setup\\_per\\_zone\\_wmarks-2.png](http://jake.dothome.co.kr/wp-content/uploads/2016/06/setup_per_zone_wmarks-2.png))

# Zone stat threaded output

High watermark to calculate stat threads for zones and nodes. The calculated items are as follows:

- node stat threaded
- John Stat Srezold
- percpu\_drift\_mark

## refresh\_zone\_stat\_thresholds()

mm/vmstat.c

```

1  | /*
2  |  * Refresh the thresholds for each zone.
3  |  */

01 | void refresh_zone_stat_thresholds(void)
02 | {
03 |     struct pglist_data *pgdat;
04 |     struct zone *zone;
05 |     int cpu;
06 |     int threshold;
07 |
08 |     /* Zero current pgdat thresholds */
09 |     for_each_online_pgdat(pgdat) {
10 |         for_each_online_cpu(cpu) {
11 |             per_cpu_ptr(pgdat->per_cpu_nodestats, cpu)->stat
12 | _threshold = 0;
13 |         }
14 |     }
15 |     for_each_populated_zone(zone) {
16 |         struct pglist_data *pgdat = zone->zone_pgdat;
17 |         unsigned long max_drift, tolerate_drift;
18 |
19 |         threshold = calculate_normal_threshold(zone);
20 |
21 |         for_each_online_cpu(cpu) {
22 |             int pgdat_threshold;
23 |
24 |             per_cpu_ptr(zone->pageset, cpu)->stat_threshold
25 |                 = threshold;
26 |
27 |             /* Base nodestat threshold on the largest popula
28 | ted zone. */
29 |             pgdat_threshold = per_cpu_ptr(pgdat->per_cpu_nod
30 | estats, cpu)->stat_thresholdd
31 | ;
32 |             per_cpu_ptr(pgdat->per_cpu_nodestats, cpu)->stat
33 | _threshold
34 |                 = max(threshold, pgdat_threshold);
35 |
36 |             /*
37 |              * Only set percpu_drift_mark if there is a danger that
38 |              * NR_FREE_PAGES reports the low watermark is ok when in
39 | fact
40 |              * the min watermark could be breached by an allocation
41 |              */
42 |             tolerate_drift = low_wmark_pages(zone) - min_wmark_pages
43 | (zone);
44 |             max_drift = num_online_cpus() * threshold;
45 |             if (max_drift > tolerate_drift)

```



```

28      * 40      16      5      900M      4
29      * 70      64      7      2-4 GB      5
30      * 84      64      7      4-8 GB      6
31      * 108     512     9      4-8 GB      6
32      * 125     1024    10     8-16 GB      8
33      * 125     1024    10     16-32 GB     9
34      */
35
36      mem = zone_managed_pages(zone) >> (27 - PAGE_SHIFT);
37
38      threshold = 2 * fls(num_online_cpus()) * (1 + fls(mem));
39
40      /*
41       * Maximum threshold is 125
42       */
43      threshold = min(125, threshold);
44
45      return threshold;
46  }

```

Proportional to the size of the managed\_pages in the request zone, the maximum number of threads to be used when under memory pressure is limited to a maximum of 125.

- Calculation Formula
  - $\text{threshold} = 2 * (\log_2(\text{number of online cpus}) + 1) * (1 + \log_2(\text{z->managed\_pages} / 128\text{M}) + 1)$ 
    - $= 2 * \text{fls}(\text{number of online CPUs}) + \text{fls}(\text{z->managed\_pages} / 128\text{M})$
  - Threshold values are limited to a maximum of 125
  - e.g. 6 CPUs, zone is managed\_pages=3.5G
    - $\text{threshold} = 2 * 3 * 5 = 30$
  - e.g. 64 cpu, zone is managed\_pages=3.5G
    - $\text{threshold} = 2 * 7 * 5 = 70$

## b) Calculating threads to be used for memory compression

The pressure thread value is the thread value used immediately after the kswapd is woken up due to insufficient memory, and the related function is as follows.

- kswapd\_try\_to\_sleep() function

### calculate\_pressure\_threshold()

mm/vmstat.c

```

01  int calculate_pressure_threshold(struct zone *zone)
02  {
03      int threshold;
04      int watermark_distance;
05
06      /*
07       * As vmstats are not up to date, there is drift between the estimated
08       * and real values. For high thresholds and a high number of CPU
09       * s, it is possible for the min watermark to be breached while the estimated
10       * value looks fine. The pressure threshold is a reduced value such

```

```

11      * that even the maximum amount of drift will not accidentally b
    reach
12      * the min watermark
13      */
14      watermark_distance = low_wmark_pages(zone) - min_wmark_pages(zon
    e);
15      threshold = max(1, (int)(watermark_distance / num_online_cpus
    ()));
16
17      /*
18       * Maximum threshold is 125
19       */
20      threshold = min(125, threshold);
21
22      return threshold;
23  }

```

Limit the maximum number of threads to be used under memory pressure in proportion to the size of the managed\_pages in the request zone, up to 125.

- In a system with a large number of CPUs, it is possible that there are a lot of PER-CPU counters that are maintained on each CPU that are not reflected in the zone counters. Therefore, in a low-memory situation, the following calculation formula is used to reduce the threaded value.
- Calculation Formula
  - threshold = interval between watermarks / number of CPUs online
  - Threshold values are limited to a maximum of 125
  - e.g. 6 CPUs, min=1975, low=2468, high=2962,
    - threshold=493 / 6 = 82
  - e.g. 64 CPUs, min=1975, low=2468, high=2962
    - threshold=493 / 64 = 7

## Lowmem Reserve Page Output

The way to specify the desired zone when requesting a memory allocation is to use the following GFP mask. If the following flags are not used, the NORMAL zone is selected by default.

- \_\_GFP\_DMA
- \_\_GFP\_DMA32
- \_\_GFP\_HIGHMEM
- \_\_GFP\_MOVABLE

Kernel developers use the following GFP mask to select zones when allocating memory:

- GFP\_DMA
  - \_\_GFP\_DMA Use the gfp flag and select the DMA zone.
- GFP\_DMA32
  - \_\_GFP\_DMA32 Use the gfp flag and select the DMA32 zone.
- GFP\_KERNEL
  - If you don't use the gfp flag associated with zone assignment, select the NORMAL zone.

- GFP\_HIGHUSER
  - \_\_GFP\_HIGHMEM Use the gfp flag and select the HIGHMEM zone.
- GFP\_HIGHUSER\_MOVABLE
  - Use the \_\_GFP\_HIGHMEM and \_\_GFP\_HIGHUSER gfp flags, and select the MOVABLE zone.

If the requested zone fails to be allocated, it will be dropped back to the next zone via the zonelist, in the following order:

- MOVABLE -> HIGHMEM -> NORMAL -> DMA32 -> DMA

## lowmem\_reserve

The memory regions used in the normal and DMA/DMA32 zones, which are lowmem zones among other zones, are pre-mapped to the kernel and allow for fast kernel memory allocation. Therefore, it is not preferred to allocate lowmem areas when requesting allocation for highmem and movable zones used by user applications. If the parent zone inevitably falls back due to insufficient memory and needs to be allocated in a zone in the lowmem zone, the number of pages to limit the allocation of that zone is specified. These values are specified as a matrix for the initial request zone and the final target zone.

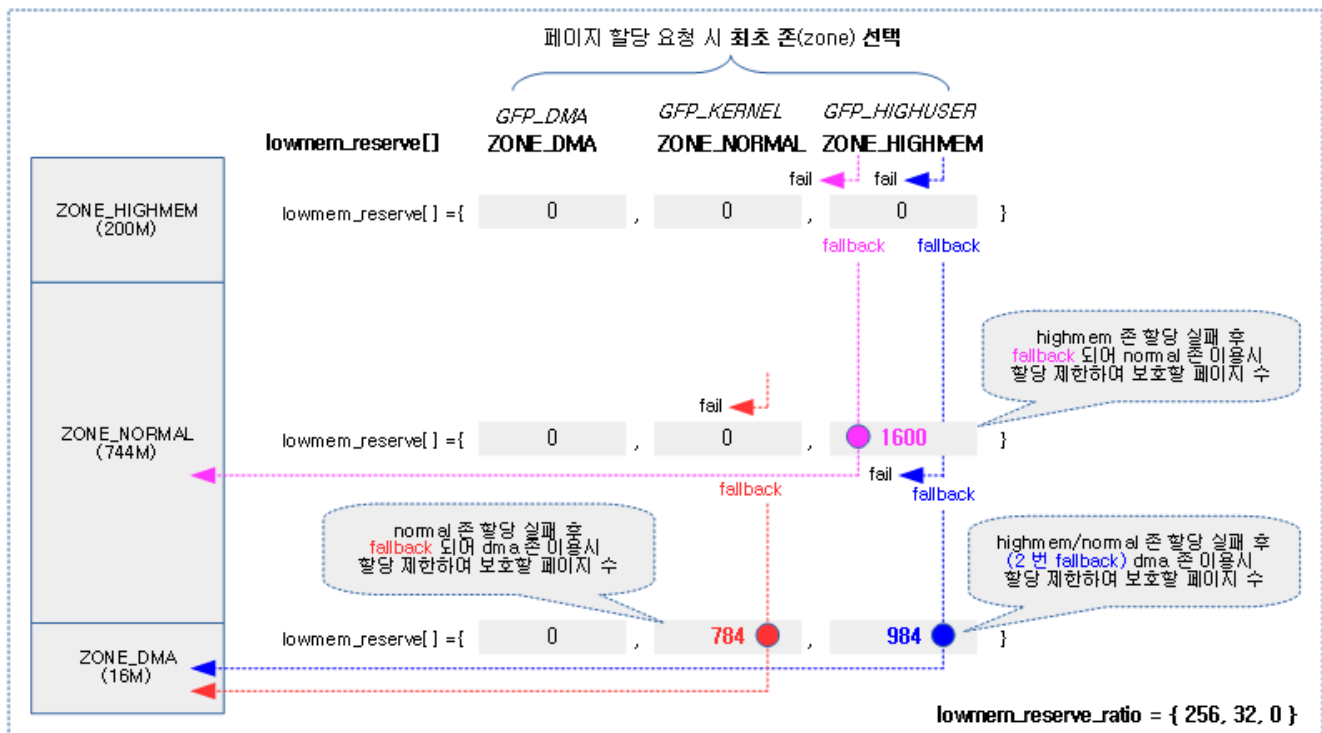
As shown in the following figure, refer to the lowmem\_reserve values of each zone.

- It is assumed that an ARM1 system with 32G memory is running three zones at dma (3M), normal (16M), and high mem (784M) respectively.
- The values in parentheses {} are the order of zones operated by the system, and those values specify the number of pages allowed to be assigned.
- If you follow the red line below, you can see that when you receive a highmem request but it fallbacks and you head to the normal zone, you can protect the lowmem area by restricting the allocation by raising it by 1600 pages higher than the watermark threshold.

최종 할당 존	최초 요청 존	ZONE_DMA	ZONE_NORMAL	ZONE_HIGHMEM	
lowmem_reserve[ZONE_DMA] =	{	0	784	984	}
lowmem_reserve[ZONE_NORMAL] =	{	0	0	1600	}
lowmem_reserve[ZONE_HIGHMEM] =	{	0	0	0	}

([http://jake.dothome.co.kr/wp-content/uploads/2016/06/lowmem\\_reserve-3.png](http://jake.dothome.co.kr/wp-content/uploads/2016/06/lowmem_reserve-3.png))

As shown in the following figure, when the allocation fails in the zone initially requested when the page is allocated, it is fallback, and when it goes to the lowmem area, it shows the page allocation limit that is added more than the watermark criterion.



([http://jake.dothome.co.kr/wp-content/uploads/2016/06/lowmem\\_reserve-1c.png](http://jake.dothome.co.kr/wp-content/uploads/2016/06/lowmem_reserve-1c.png))

## lowmem\_reserve\_ratio

It is a ratio for calculating the `lowmem_reserve` value for each zone using the `managed_pages` of each zone.

You can check the `lowmem_reserve_ratio` values as follows:

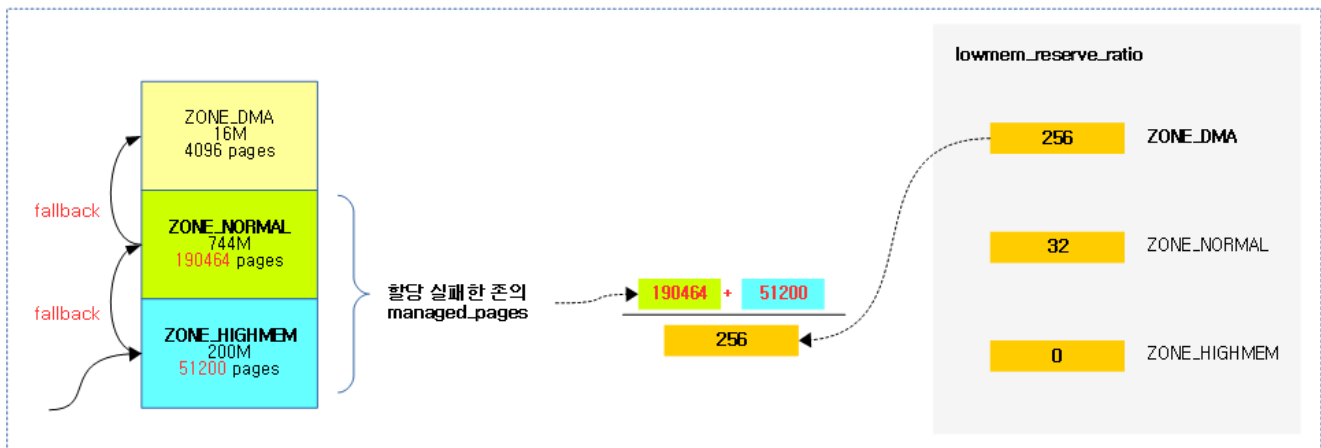
- The ratio is set in order from the lowest to the highest zone operated in the system.
- The example values below are in the order of DMA, Normal, and Highmem zones.

```
1 | $ cat /proc/sys/vm/lowmem_reserve_ratio
2 | 256      32      0
```

As shown in the following figure, the calculation of the number of pages allowed when fallback for the allocation request zone is as follows:

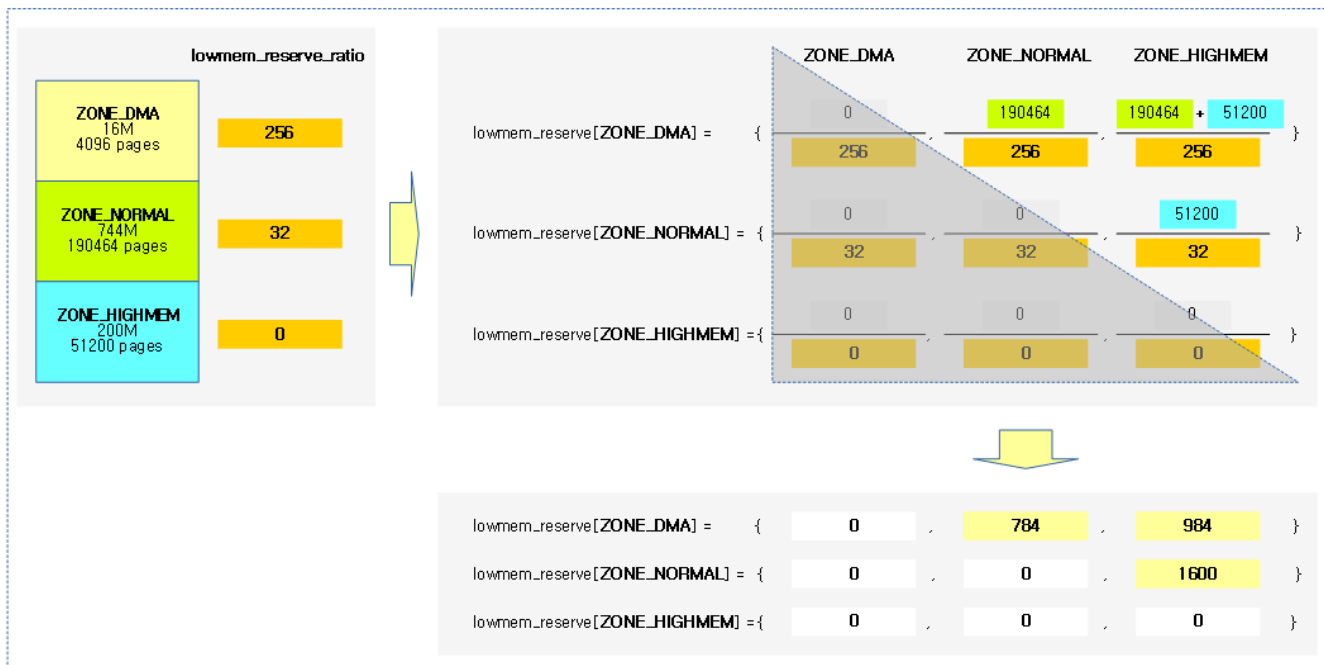
- It can be seen that the value of the numerator adds the number of managed\_pages from the request zone to the allocation failure zone.
- It can be seen that it is fallback to the value of the denominator, specifying the lowmem\_reserve\_ratio of the final allocation zone.





([http://jake.dothome.co.kr/wp-content/uploads/2016/06/lowmem\\_reserve-4.png](http://jake.dothome.co.kr/wp-content/uploads/2016/06/lowmem_reserve-4.png))

As shown in the following figure, the `managed_pages` size and `lowmem_reserve_ratio` of each zone are used to calculate the `lowmem_reserve` value of each zone.



([http://jake.dothome.co.kr/wp-content/uploads/2016/06/lowmem\\_reserve-2.png](http://jake.dothome.co.kr/wp-content/uploads/2016/06/lowmem_reserve-2.png))

Next, let's look at the node-zone `lowmem_reserve` values on an x4 system with NUMA and four zones (dma, dma32, normal, movable).

- Take a look at the protection section below. When this feature was first developed, the variable was called `protection`, which was later changed to `lowmem_reserve`.

```
01 $ cat /proc/zoneinfo
02 Node 0, zone    DMA
03   pages free    3915
04   min          5
05   low          6
06   high         7
07   scanned      0
08   spanned      4095
09   present      3992
10   managed      3971
11   nr_free_pages 3915
```

```

12 ...
13     protection: (0, 1675, 31880, 31880)
14 ...
15 Node 0, zone    DMA32
16   pages free    59594
17     min      588
18     low      735
19     high     882
20     scanned   0
21     spanned  1044480
22     present   491295
23     managed   429342
24     nr_free_pages 59594
25 ...
26     protection: (0, 0, 30204, 30204)
27 ...
28 Node 0, zone    Normal
29   pages free    902456
30     min     10607
31     low     13258
32     high    15910
33     scanned   0
34     spanned  7864320
35     present  7864320
36     managed  7732469
37     nr_free_pages 902456
38 ...
39     protection: (0, 0, 0, 0)
40 ...
41 Node 1, zone    Normal
42   pages free   1133093
43     min     11326
44     low     14157
45     high    16989
46     scanned   0
47     spanned  8388608
48     present  8388608
49     managed  8256697
50     nr_free_pages 1133093
51 ...
52     protection: (0, 0, 0, 0)
53 ...

```

You can see that the `lowmem_reserve` value in a system with only one zone is 0 as follows.

- In the case of an ARM32 system that operates the `dma64` and `normal` zones, it operates only one zone if all the memory is used with only the `dma32` zone. In other words, there is no fallback.

```

$ cat /proc/zoneinfo
Node 0, zone    DMA32
  per-node stats
  (...생략...)
    pages free      633680
      min          5632
      low          7040
      high         8448
    spanned 786432
    present 786432
    managed 765771
    protection: (0, 0, 0)    <-----
    nr_free_pages 633680
  (...생략...)
Node 0, zone    Normal
  pages free      0
    min          0
    low          0
    high         0
  spanned 0
  present 0
  managed 0
  protection: (0, 0, 0)    <-----
  (...생략...)

```

## setup\_per\_zone\_lowmem\_reserve()

mm/page\_alloc.c

```

1  /*
2  * setup_per_zone_lowmem_reserve - called whenever
3  *   sysctl_lower_zone_reserve_ratio changes. Ensures that each zone
4  *   has a correct pages reserved value, so an adequate number of
5  *   pages are left in the zone after a successful __alloc_pages().
6  */

01 static void setup_per_zone_lowmem_reserve(void)
02 {
03     struct pglist_data *pgdat;
04     enum zone_type j, idx;
05
06     for_each_online_pgdat(pgdat) {
07         for (j = 0; j < MAX_NR_ZONES; j++) {
08             struct zone *zone = pgdat->node_zones + j;
09             unsigned long managed_pages = zone_managed_pages
10 (zone);
11
12             zone->lowmem_reserve[j] = 0;
13
14             idx = j;
15             while (idx) {
16                 struct zone *lower_zone;
17
18                 idx--;
19                 lower_zone = pgdat->node_zones + idx;
20
21                 if (sysctl_lowmem_reserve_ratio[idx] <
22 = 0;

```

```

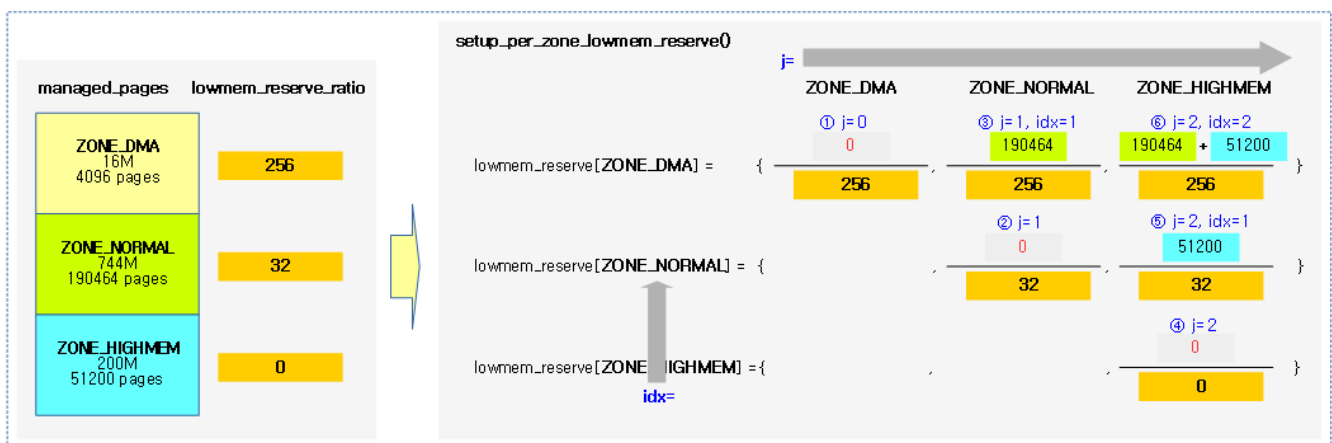
22     lower_zone->lowmem_reserve[j] =
23     0;
24     } else {
25         lower_zone->lowmem_reserve[j] =
26             managed_pages / sysctl_l
27             owmem_reserve_ratio[idx];
28     }
29     managed_pages += zone_managed_pages(love
30     r_zone);
31     }
32     }
33     /* update totalreserve_pages */
34     calculate_totalreserve_pages();

```

The managed\_pages value of the zone and the ratio of the lowmem\_reserve\_ratio are used to calculate the lowmem\_reserve[] page values for each zone.

- In code lines 6~11, traverse the total number of online nodes, internally traverse the number of available zones in that node, and reset the lowmem\_reserve[] value to 0.
- In code lines 13~30, go around the loop from the zone designated as the loop counter except for the first zone, and add the managed\_pages except the first zone (idx=0) to specify the lowmem\_reserve\_ratio ratio to the lowmem\_reserve.
- Since the lowmem\_reserve value has changed in line 33 of code, the number of TotalReserve pages is also updated.

The following figure shows the process of calculating the lowmem\_reserve[] for each zone.



([http://jake.dothome.co.kr/wp-content/uploads/2016/06/setup\\_per\\_zone\\_lowmem\\_reserve-2.png](http://jake.dothome.co.kr/wp-content/uploads/2016/06/setup_per_zone_lowmem_reserve-2.png))

## TotalReserve Page Output

### calculate\_totalreserve\_pages()

mm/page\_alloc.c

```

1  /*
2  * calculate_totalreserve_pages - called when sysctl_lower_zone_reserve_
3  * ratio
4  * or min_free_kbytes changes.

```

4 | \*/

```

01 static void calculate_totalreserve_pages(void)
02 {
03     struct pglist_data *pgdat;
04     unsigned long reserve_pages = 0;
05     enum zone_type i, j;
06
07     for_each_online_pgdat(pgdat) {
08
09         pgdat->totalreserve_pages = 0;
10
11         for (i = 0; i < MAX_NR_ZONES; i++) {
12             struct zone *zone = pgdat->node_zones + i;
13             long max = 0;
14             unsigned long managed_pages = zone_managed_pages
(zone);
15
16             /* Find valid and maximum lowmem_reserve in the
zone */
17             for (j = i; j < MAX_NR_ZONES; j++) {
18                 if (zone->lowmem_reserve[j] > max)
19                     max = zone->lowmem_reserve[j];
20             }
21
22             /* we treat the high watermark as reserved page
s. */
23             max += high_wmark_pages(zone);
24
25             if (max > managed_pages)
26                 max = managed_pages;
27
28             pgdat->totalreserve_pages += max;
29
30             reserve_pages += max;
31         }
32     }
33     totalreserve_pages = reserve_pages;
34 }

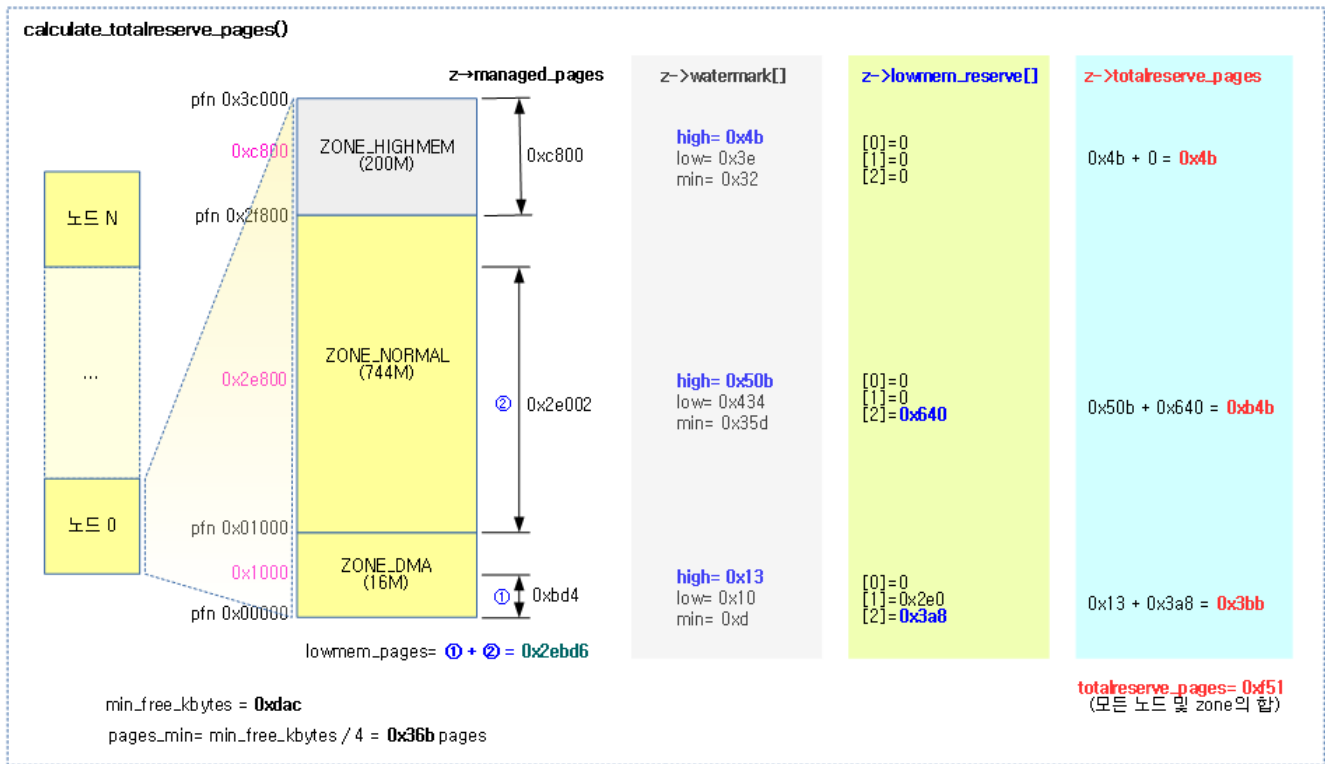
```

Calculate the totalreserve\_pages value for each node and the global totalreserve\_pages value.

- In line 7~9 of the code, it traverses all nodes and initializes them to 0 first to calculate the totalreserve\_pages.
- In lines 11~14 of code, traverse all zones and find out the managed page of that zone.
- From code lines 17~28, find out the maximum value of the lowmem\_reserve set from the traversing zone to the last zone, and then add a high watermark to add it to the totalreserve\_pages. Limit the value to be added so that it does not exceed the managed\_pages of the zone in the circuit.
- In line 30~33 of the code, assign the values added by each node and zone to the global totalreserve\_pages.

The following figure shows the process of calculating the following values by adding each node and zone to calculate the totalreserve\_pages value.

- Maximum lowmem\_reserve above each zone
- High watermark for each zone



([http://jake.dothome.co.kr/wp-content/uploads/2016/06/calculate\\_totalreserve\\_pages-1b.png](http://jake.dothome.co.kr/wp-content/uploads/2016/06/calculate_totalreserve_pages-1b.png))

## inactive\_ratio Output

The following steps removed the `inactive_ratio` from kernel v4.7-rc1 on bootup. The proportion of inactive lists is now deprecated because it is automatically compared to the entire memory in the `inactive_list_is_low()` function.

- 참고: mm: vmscan: reduce size of inactive file list  
<https://github.com/torvalds/linux/commit/59dc76b0d4dfdd7dc46a1010e4afb44f60f3e97f#diff-d1793bdc2ce9e21814061bb882d0c3ac>

## setup\_per\_zone\_inactive\_ratio()

mm/page\_alloc.c

```

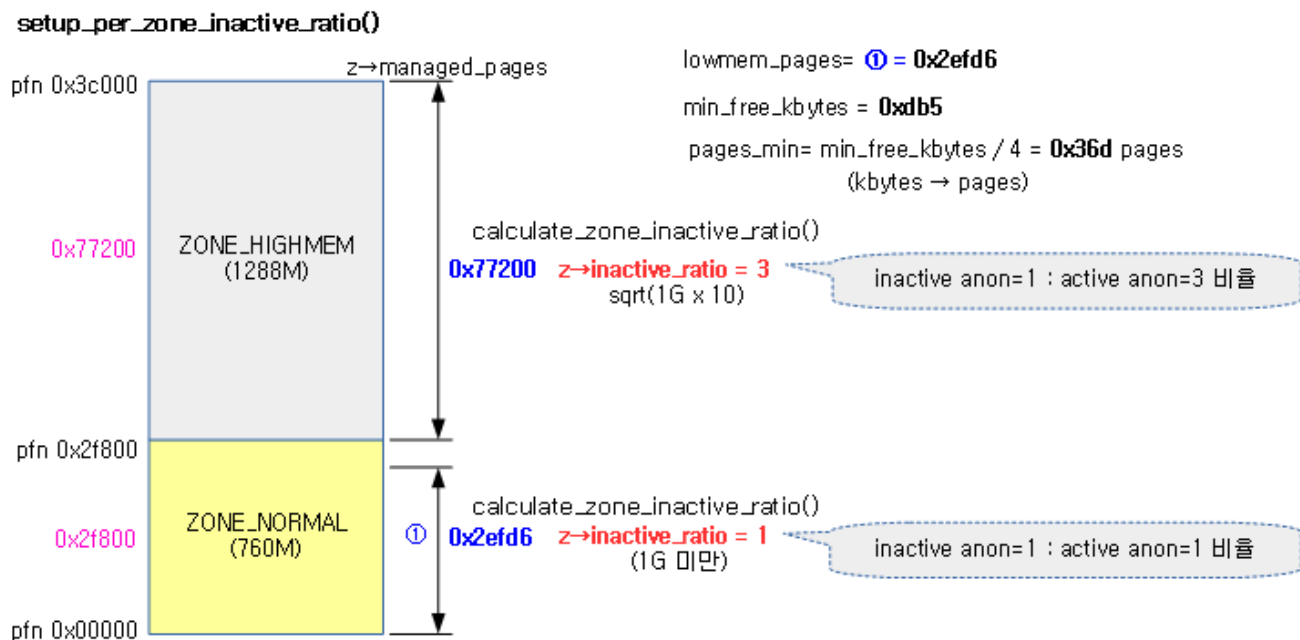
1 | static void __meminit setup_per_zone_inactive_ratio(void)
2 | {
3 |     struct zone *zone;
4 |     for_each_zone(zone)
5 |         calculate_zone_inactive_ratio(zone);
6 | }

```

Calculate the inactive anon rate of all zones.

- If the `zone->managed_pages` is less than 256k pages (1GB), the `zone->inactive_ratio` will be 1 and the ratio of active and inactive anon will be set to 1:1. If the `managed_pages` exceeds 256k pages (1GB), the `inactive_ratio` value will be greater than 3 and the percentage of inactive anon will be reduced to less than 1/3.

The following figure shows the process of calculating inactive\_ratio values based on the managed\_pages size of each zone.



([http://jake.dothome.co.kr/wp-content/uploads/2016/06/setup\\_per\\_zone\\_inactive\\_ratio-1.png](http://jake.dothome.co.kr/wp-content/uploads/2016/06/setup_per_zone_inactive_ratio-1.png))

## calculate\_zone\_inactive\_ratio()

mm/page\_alloc.c

```

01  /*
02  * The inactive anon list should be small enough that the VM never has t
03  * do too much work, but large enough that each inactive page has a chan
04  * ce
05  * to be referenced again before it is swapped out.
06  *
07  * The inactive_anon ratio is the target ratio of ACTIVE_ANON to
08  * INACTIVE_ANON pages on this zone's LRU, maintained by the
09  * pageout code. A zone->inactive_ratio of 3 means 3:1 or 25% of
10  * the anonymous pages are kept on the inactive list.
11  *
12  * total      target      max
13  * memory     ratio      inactive anon
14  * -----
15  * 10MB       1          5MB
16  * 100MB      1          50MB
17  * 1GB        3          250MB
18  * 10GB       10         0.9GB
19  * 100GB      31         3GB
20  * 1TB        101        10GB
21  * 10TB       320        32GB
22  */
23  static void __meminit calculate_zone_inactive_ratio(struct zone *zone)
24  {
25      unsigned int gb, ratio;
26
27      /* Zone size in gigabytes */
28      gb = zone->managed_pages >> (30 - PAGE_SHIFT);
29      if (gb)
30          ratio = int_sqrt(10 * gb);
31      else

```

```

31         ratio = 1;
32
33     zone->inactive_ratio = ratio;
34 }

```

Calculate the inactive anon rate of the specified zone.

- If the managed page is less than 1G, the ratio will be 1.
  - inactive anon=1 : active anon=1
- If the managed page is more than 1G, set the ratio by using \* 10 and then the square root (root).  
(inactive = 1 / ratio)
  - 예) managed pages=1G
    - inactive\_ratio=3
      - inactive anon=1 : active anon=3
      - inactive anon=250M : active anon=750M

## consultation

- Zoned Allocator -1- (Physics Page Assignment - Fastpath) (<http://jake.dothome.co.kr/zoned-allocator-alloc-pages-fastpath>) | Qc
- Zoned Allocator -2- (Physics Page Assignment - Slowpath) (<http://jake.dothome.co.kr/zoned-allocator-alloc-pages-slowpath>) | Qc
- Zoned Allocator -3- (Buddy Page Allocation) (<http://jake.dothome.co.kr/buddy-alloc>) | Qc
- Zoned Allocator -4- (Buddy Page Terminated) (<http://jake.dothome.co.kr/buddy-free/>) | Qc
- Zoned Allocator -5- (Per-CPU Page Frame Cache) (<http://jake.dothome.co.kr/per-cpu-page-frame-cache>) | 문c
- Zoned Allocator -6- (Watermark) (<http://jake.dothome.co.kr/zoned-allocator-watermark>) | Sentence C – Current post
- Zoned Allocator -7- (Direct Compact) (<http://jake.dothome.co.kr/zoned-allocator-compaction>) | 문c
- Zoned Allocator -8- (Direct Compact-Isolation) (<http://jake.dothome.co.kr/zoned-allocator-isolation>) | 문c
- Zoned Allocator -9- (Direct Compact-Migration) (<http://jake.dothome.co.kr/zoned-allocator-migration>) | 문c
- Zoned Allocator -10- (LRU & pagevec) (<http://jake.dothome.co.kr/lru-lists-pagevecs>) | 문c
- Zoned Allocator -11- (Direct Reclaim) (<http://jake.dothome.co.kr/zoned-allocator-reclaim>) | 문c
- Zoned Allocator -12- (Direct Reclaim-Shrink-1) (<http://jake.dothome.co.kr/zoned-allocator-shrink-1>) | 문c
- Zoned Allocator -13- (Direct Reclaim-Shrink-2) (<http://jake.dothome.co.kr/zoned-allocator-shrink-2>) | 문c
- Zoned Allocator -14- (Kswapd) (<http://jake.dothome.co.kr/zoned-allocator-kswapd>) | 문c



# 11 thoughts to “Zoned Allocator -6- (Watermark)”



**AUSTIN KIM (HTTP://ROUSALOME.EGLOOS.COM/)**

2020-01-13 14:18 (<http://jake.dothome.co.kr/zoned-allocator-watermark/#comment-228690>)

This is Moon Young-il and Kim Dong-hyun.

I'm studying and learning a lot from the blog. Thank you as always.

However, I saw a typo in line 3 of the WATERMARK\_MIN entry. I would like to contribute to the 'Munc Blog'.

'Page recall price -> page reclaim'

RESPONSE (/ZONNED-ALLOCATOR-WATERMARK/?REPLYTOCOM=228690#RESPOND)



**MOON YOUNG-IL (HTTP://JAKE.DOTHOME.CO.KR)**

2020-01-15 07:53 (<http://jake.dothome.co.kr/zoned-allocator-watermark/#comment-228903>)

Donghyun, how did you spend the end of the year?

In my first job in the 1990s, I wondered, 'Will the year 2000 come?'

But every day, I see the number 2000, which I thought would never come, let alone 2020.

It's a big new year with a two-digit change, so I wish Donghyun a Happy New Year as well.

I'm sure you're busy with the publication of your kernel book, and I hope your life will be richer this year.

We hope to see you often. ^^

RESPONSE (/ZONNED-ALLOCATOR-WATERMARK/?REPLYTOCOM=228903#RESPOND)



**AUSTIN KIM (HTTP://ROUSALOME.EGLOOS.COM)**

2020-01-16 21:44 (<http://jake.dothome.co.kr/zoned-allocator-watermark/#comment-228987>)

Thank you for your support.

It's late, but I wish you a happy new year and I hope to see you again this year with a different opportunity.

Like the 'Munc Blog', I will study kernel diligently this year.

RESPONSE (/ZONNED-ALLOCATOR-WATERMARK/?REPLYTOCOM=228987#RESPOND)



**PETER**

2020-12-13 13:48 (<http://jake.dothome.co.kr/zoned-allocator-watermark/#comment-294600>)

I have a question in the `__zone_watermark_ok()` function.

"Reduce the number of pages requested to be allocated free pages on line 10 by -1."

If you count the number of pages with the requested order, subtract 1 and put it in the `free_pages`,

why do you do -1 here?

응답 (/ZONNED-ALLOCATOR-WATERMARK/?REPLYTOCOM=294600#RESPOND)



**문영일 (HTTP://JAKE.DOTHOME.CO.KR)**

2020-12-13 14:54 (<http://jake.dothome.co.kr/zoned-allocator-watermark/#comment-294609>)

안녕하세요?

저도 예전에 한 번 왜 -1을 할 까 생각을 한적이 있었습시다만

개발자의 정확한 의미는 모르고 지나간 적이 있었습시다. 다만 아래 코드와 같이 코드를 작성해도 되는데

```
free_pages -= (1 < < order); ... if ((free_pages + 1) <= min + z-  
>lowmem_reserve[classzone_idx])
```

커널 코드는 다음과 같이 미리 1을 덜 빼두고 비교 판단 코드를 작성하기 위함이라 판단됩니다.

감사합니다.

```
free_pages -= (1 < < order) - 1; ... if (free_pages <= min + z-  
>lowmem_reserve[classzone_idx])
```

응답 (/ZONNED-ALLOCATOR-WATERMARK/?REPLYTOCOM=294609#RESPOND)



**PETER**

2020-12-13 19:44 (<http://jake.dothome.co.kr/zoned-allocator-watermark/#comment-294662>)

생각해 보니 아래와 같이 등호(=)를 빼는 방법도 가능해 보입니다.

```
free_pages -= (1 < < order); ... if ((free_pages) lowmem_reserve[classzone_idx])
```

말씀하신 위치의 비교 조건에만 관련된 구현의 차이 일뿐,

알고리즘상의 큰 차이는 없어 보입니다.

감사합니다.

응답 (/ZONNED-ALLOCATOR-WATERMARK/?REPLYTOCOM=294662#RESPOND)

**PETER**2020-12-13 19:46 (<http://jake.dothome.co.kr/zoned-allocator-watermark/#comment-294663>)

if (free\_pages lowmem\_reserve[classzone\_idx])  
입니다

응답 (/ZONNED-ALLOCATOR-WATERMARK/?REPLYTOCOM=294663#RESPOND)

**문영일 (HTTP://JAKE.DOTHOME.CO.KR)**2020-12-14 06:36 (<http://jake.dothome.co.kr/zoned-allocator-watermark/#comment-294778>)

네. 부등호(=)를 빼도 될 듯 합니다. ^^

응답 (/ZONNED-ALLOCATOR-WATERMARK/?REPLYTOCOM=294778#RESPOND)

**메모리바보**2021-01-22 22:35 (<http://jake.dothome.co.kr/zoned-allocator-watermark/#comment-303303>)

안녕하세요,

좋은 정보 감사합니다.

본문에 나오는 lowmem은 zone\_highmem을 제외한 영역인가요? 그렇다면 lowmem을 기준으로 watermark를 산정하는 이유는 무엇인가요?

응답 (/ZONNED-ALLOCATOR-WATERMARK/?REPLYTOCOM=303303#RESPOND)

**문영일 (HTTP://JAKE.DOTHOME.CO.KR)**2021-01-23 14:09 (<http://jake.dothome.co.kr/zoned-allocator-watermark/#comment-303487>)

안녕하세요?

커널에서 사용되는 lowmem은 다음 두 가지 의미로 사용됩니다.

- 1) 메모리 존을 뜻할 때에는 highmem 존의 반대 의미로 사용되고, lowmem에는 normal, dma32, dma 존이 포함됩니다. (64bit의 경우 highmem이 존재하지 않습니다)
- 2) 메모리 부족 경계를 의미로 사용되는 lowmem도 있습니다. 따라서 이 기준에 워터마크를 사용합니다.

감사합니다.

응답 (/ZONNED-ALLOCATOR-WATERMARK/?REPLYTOCOM=303487#RESPOND)

**메모리바보**2021-01-23 19:24 (<http://jake.dothome.co.kr/zonned-allocator-watermark/#comment-303536>)

아 두가지가 다른 의미를 가지는군요.

좋은 정보 감사합니다.!

응답 (/ZONNED-ALLOCATOR-WATERMARK/?REPLYTOCOM=303536#RESPOND)

**댓글 남기기**

이메일은 공개되지 않습니다. 필수 입력창은 \* 로 표시되어 있습니다

댓글

이름 \*

이메일 \*

웹사이트

댓글 작성

[◀ Zoned Allocator -2- \(물리 페이지 할당-Slowpath\) \(http://jake.dothome.co.kr/zonned-allocator-alloc-pages-slowpath/\)](http://jake.dothome.co.kr/zonned-allocator-alloc-pages-slowpath/)[Zoned Allocator -7- \(Direct Compact\) ▶ \(http://jake.dothome.co.kr/zonned-allocator-compaction/\)](http://jake.dothome.co.kr/zonned-allocator-compaction/)

문c 블로그 (2015 ~ 2023)