

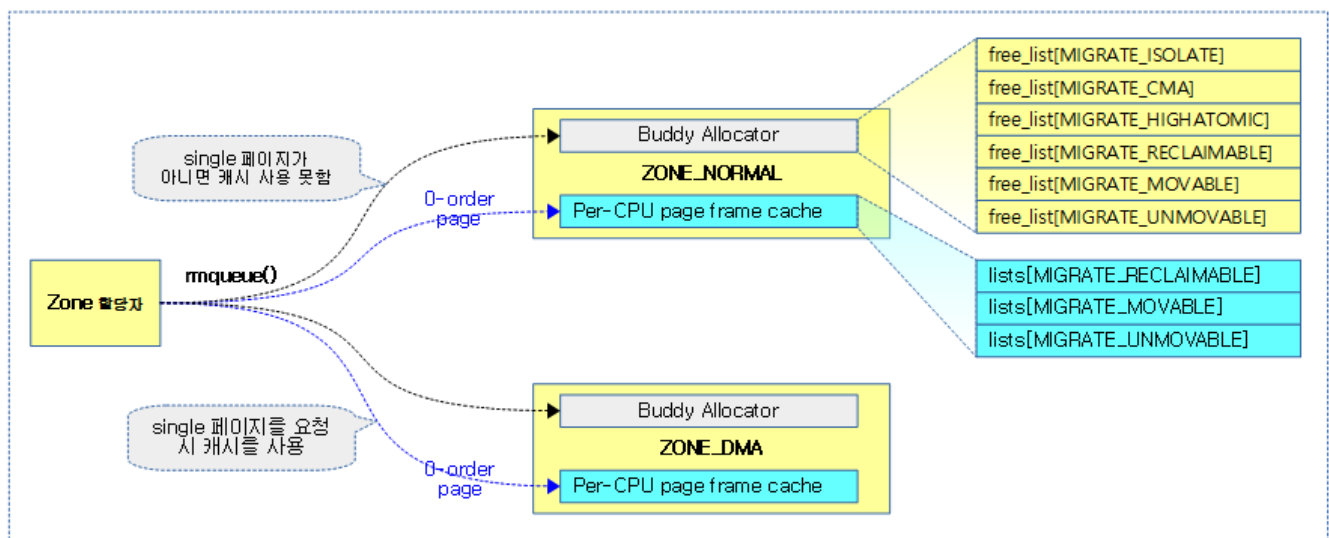
Zoned Allocator -5- (Per-CPU Page Frame Cache)

📅 2016-04-21 (<http://jake.dothome.co.kr/per-cpu-page-frame-cache/>) 👤 Moon Young-il

(<http://jake.dothome.co.kr/author/admin/>) 📁 Linux Kernel (<http://jake.dothome.co.kr/category/linux/>)

<kernel v5.0>

Per-CPU Page Frame Cache



(<http://jake.dothome.co.kr/wp-content/uploads/2016/04/per-cpu-page-frame-cache-1d.png>)

- In the kernel, memory allocation is usually a request for a single page (0-order page) frame rather than a large page.
- In order to increase the allocation processing performance only in the case of a single page request, a per-CPU page frame cache is prepared in each zone, multiple pages are prepared in advance, and the cached page is delivered directly to the requester without using buddy for the requested single page.
 - When using Buddy, it is necessary to lock the zone, but it was implemented as lock-less using Per-CPU to improve performance.
- In the previous kernel, each zone used to use a hot cache and a cold cache, but the purpose has been consolidated into one, and the front of the prepared cache page is used for hot requests, and the back page is used for cold requests.
- Since kernel 2.6.25-rc1, each zone is managed as a cache array with a number of 3 migratetypes.
- If the kernel monitor detects that the cache page falls below the low watermark, it allocates a batch of pages to the cache in advance.
- If there are no cached pages in a single page request, batch of pages are allocated to the cache and then one of them is delivered to the requester.
- If the cached page is higher than high when a single page is released, it returns the number of batches back to the buddy system.

Allocate/unallocate/unassign order 0 pages in PCP

Order 0 Page Assignment

rmqueue_pcplist()

mm/page_alloc.c

```

01  /* Lock and remove page from the per-cpu list */
02  static struct page *rmqueue_pcplist(struct zone *preferred_zone,
03                                     struct zone *zone, unsigned int order,
04                                     gfp_t gfp_flags, int migratetype,
05                                     unsigned int alloc_flags)
06  {
07      struct per_cpu_pages *pcp;
08      struct list_head *list;
09      struct page *page;
10      unsigned long flags;
11
12      local_irq_save(flags);
13      pcp = &this_cpu_ptr(zone->pageset)->pcp;
14      list = &pcp->lists[migratetype];
15      page = __rmqueue_pcplist(zone, migratetype, alloc_flags, pcp, l
16  ist);
17      if (page) {
18          __count_zid_vm_events(PGALLOC, page_zonenum(page), 1 <<
19  order);
20          zone_statistics(preferred_zone, zone);
21      }
22      local_irq_restore(flags);
23      return page;
24  }

```

@migratetype assigns page 0 of order <> from pcp and returns a page descriptor.

- PCP, which is a buddy system cache implemented with Per-CPU in line 12~15 of code, allocates the order 0 page of the @migratetype from PCP with only local interrupts disabled without using a cost-intensive lock.
- In lines 16~19 of code, increment the PGALLOC counter by the number of pages.
- Returns the page assigned in line 21 of code.

__rmqueue_pcplist()

mm/page_alloc.c

```

01  /* Remove page from the per-cpu list, caller must protect the list */
02  static struct page *__rmqueue_pcplist(struct zone *zone, int migratetype,
03  e,
04  unsigned int alloc_flags,
05  struct per_cpu_pages *pcp,
06  struct list_head *list)
07  {
08      struct page *page;
09
10      do {
11          if (list_empty(list)) {
12              pcp->count += rmqueue_bulk(zone, 0,
13              pcp->batch, list,
14              migratetype, alloc_flags);
15          }
16          if (unlikely(list_empty(list)))
17              continue;
18          page = list_first_entry(list, struct page, lru);
19          list_del(page);
20          return page;
21      } while (1);
22  }

```

```

15         return NULL;
16     }
17
18     page = list_first_entry(list, struct page, lru);
19     list_del(&page->lru);
20     pcp->count--;
21     } while (check_new_pcp(page));
22
23     return page;
24 }

```

@migratetype allocates order 0 pages from pcp and returns a page descriptor.

- If the @list of pcp is empty in line 9~16, migrate the number of pcp->batches from the buddy system.
- In code lines 18~20, get the first entry from the @list of PCP.
- Check that there are no problems with the entry to be assigned at line 21 of the code.
- Returns page 23 of order to be assigned in line 0 of code.

Order Recall 0 Pages

free_unref_page()

mm/page_alloc.c

```

1  /*
2   * Free a 0-order page
3   */
01 void free_unref_page(struct page *page)
02 {
03     unsigned long flags;
04     unsigned long pfn = page_to_pfn(page);
05
06     if (!free_unref_page_prepare(page, pfn))
07         return;
08
09     local_irq_save(flags);
10     free_unref_page_commit(page, pfn);
11     local_irq_restore(flags);
12 }

```

Order 0 page back to pcp.

- Prepare the page to be free in line 6~7 of the code. If the page status is judged to be bad, the function exits.
- In line 9~11 of code, the order 0 page is returned to PCP with the local irq disabled.

free_unref_page_prepare()

mm/page_alloc.c

```

01 static bool free_unref_page_prepare(struct page *page, unsigned long pfn)
02 {
03     int migratetype;
04
05     if (!free_pcp_prepare(page))

```

```

06         return false;
07
08         migratetype = get_pfnblock_migratetype(page, pfn);
09         set_pcpage_migratetype(page, migratetype);
10         return true;
11     }

```

Prepare a page to free. (normal=true, bad=false)

- Check the status of the page to be free in line 5~6 of the code, and if it is judged bad, it will return a false result.
- In lines 8~10 of the code, store the migrate type of the page block to which the page belongs to the page and return true.

free_unref_page_commit()

mm/page_alloc.c

```

01 static void free_unref_page_commit(struct page *page, unsigned long pfn)
02 {
03     struct zone *zone = page_zone(page);
04     struct per_cpu_pages *pcp;
05     int migratetype;
06
07     migratetype = get_pcpage_migratetype(page);
08     __count_vm_event(PGFREE);
09
10     /*
11      * We only track unmovable, reclaimable and movable on pcp list
12      * Free ISOLATE pages back to the allocator because they are bei
13      * ng
14      * offlined but treat HIGHATOMIC as movable pages so we can get
15      * those
16      * areas back if necessary. Otherwise, we may have to free
17      * excessively into the page allocator
18      */
19     if (migratetype >= MIGRATE_PCPTYPES) {
20         if (unlikely(is_migrate_isolate(migratetype))) {
21             free_one_page(zone, page, pfn, 0, migratetype);
22             return;
23         }
24         migratetype = MIGRATE_MOVABLE;
25
26     pcp = &this_cpu_ptr(zone->pageset)->pcp;
27     list_add(&page->lru, &pcp->lists[migratetype]);
28     pcp->count++;
29     if (pcp->count >= pcp->high) {
30         unsigned long batch = READ_ONCE(pcp->batch);
31         free_pcpages_bulk(zone, batch, pcp);
32     }

```

Reclaim the 0-order page to be free to pcp.

- Increment the PGFREE counter on line 8 of code.
- In code lines 17~23, the isolate type is recalled to the buddy system, and the remaining CMA and highatomic types that are not handled by PCP are changed to movable types.
- In line 25~27 of code, add it to the pcp of type migrate.

- In line 28~31 of the code, the number of entries in the PCP list is PCP->high or higher. In order to manage only a certain amount of PCP, the number of PCP->batches is transferred to the buddy system.

get_pcpage_migratetype()

include/linux/mm.h

```

1  /*
2   * A cached value of the page's pageblock's migratetype, used when the p
3   * age is
4   * put on a pcplist. Used to avoid the pageblock migratetype lookup when
5   * freeing from pcplists in most cases, at the cost of possibly becoming
6   * stale.
7   * Also the migratetype set in the page does not necessarily match the p
8   * cplist
9   * index, e.g. page might have MIGRATE_CMA set but be on a pcplist with
10  * any
11  * other index - this ensures that it will be put on the correct CMA fre
12  * elist.
13  */
14
15 static inline int get_pcpage_migratetype(struct page *page)
16 {
17     return page->index;
18 }

```

Finds the migratetype stored in page->index.

PCP <-> Buddy System Bulk Allocation/Retrieval

PCP <-Buddy System Bulk Allocation

rmqueue_bulk()

mm/page_alloc.c

```

1  /*
2   * Obtain a specified number of elements from the buddy allocator, all u
3   * nder
4   * a single hold of the lock, for efficiency. Add them to the supplied
5   * list.
6   * Returns the number of new pages which were placed at *list.
7   */
8
9 static int rmqueue_bulk(struct zone *zone, unsigned int order,
10                        unsigned long count, struct list_head *list,
11                        int migratetype, unsigned int alloc_flags)
12 {
13     int i, allocated = 0;
14
15     spin_lock(&zone->lock);
16     for (i = 0; i < count; ++i) {
17         struct page *page = __rmqueue(zone, order, migratetype,
18                                       alloc_flags);
19         if (unlikely(page == NULL))
20             break;
21     }
22 }

```

```

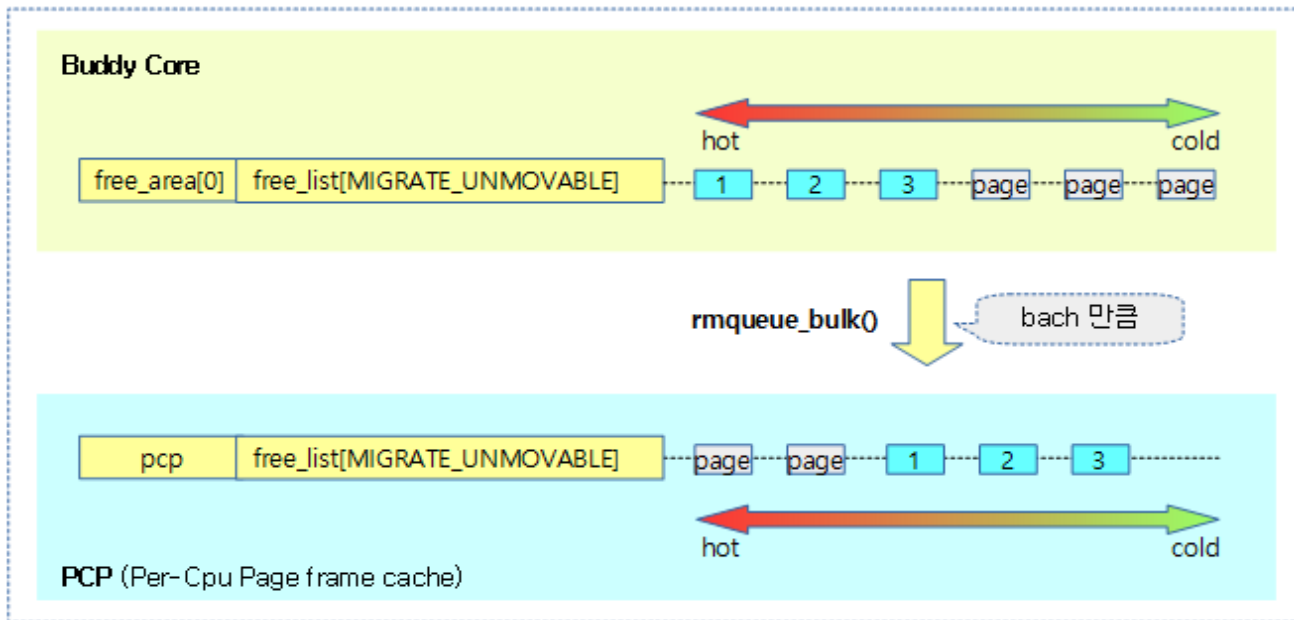
14         if (unlikely(check_pcp_refill(page)))
15             continue;
16
17         /*
18          * Split buddy pages returned by expand() are received here in
19          * physical page order. The page is added to the tail of
20          * caller's list. From the callers perspective, the link
21          * is ordered by page number under some conditions. This
22          * is useful for IO devices that can forward direction from
23          * the head, thus also in the physical page order. This is useful
24          * for IO devices that can merge IO requests if the physical
25          * pages are ordered properly.
26          */
27         list_add_tail(&page->lru, list);
28         allocated++;
29         if (is_migrate_cma(get_pcppage_migratetype(page)))
30             __mod_zone_page_state(zone, NR_FREE_CMA_PAGES,
31                                   -(1 << order));
32     }
33
34     /*
35      * i pages were removed from the buddy list even if some leak due
36      * to check_pcp_refill failing so adjust NR_FREE_PAGES based
37      * on i. Do not confuse with 'allocated' which is the number of
38      * pages added to the pcp list.
39      */
40     __mod_zone_page_state(zone, NR_FREE_PAGES, -(i << order));
41     spin_unlock(&zone->lock);
42     return allocated;
43 }

```

The Buddy system takes @count free pages from the @order slot and moves them to the @list. It then returns the actual number moved.

- Loop around the number of counts in lines 8~15 to get @order pages from the buddy system.
- Add the page from lines 27~28 to the @list.
- In code lines 29~31, if it is a cma page, decrement the NR_FREE_CMA_PAGES counter by the number of pages.
- At line 40 of code, decrement the NR_FREE_PAGES by the number of pages moved in the loop.
- Returns the number shifted in line 42 of code.

The following figure shows the bulk movement of free pages in the buddy system to PCP in batches.



(http://jake.dothome.co.kr/wp-content/uploads/2016/04/rmqueue_bulk-1b.png)

PCP -> Buddy System Bulk Recovery

free_pcppages_bulk()

```

01  /*
02  * Frees a number of pages from the PCP lists
03  * Assumes all pages on list are in same zone, and of same order.
04  * count is the number of pages to free.
05  *
06  * If the zone was previously in an "all pages pinned" state then look t
07  * see if this freeing clears that state.
08  *
09  * And clear the zone's pages_scanned counter, to hold off the "all page
10  * s are
11  * pinned" detection logic.
12  */
13
14  static void free_pcppages_bulk(struct zone *zone, int count,
15                                struct per_cpu_pages *pcp)
16  {
17      int migratetype = 0;
18      int batch_free = 0;
19      int prefetch_nr = 0;
20      bool isolated_pageblocks;
21      struct page *page, *tmp;
22      LIST_HEAD(head);
23
24      while (count) {
25          struct list_head *list;
26
27          /*
28           * Remove pages from lists in a round-robin fashion. A
29           * batch_free count is maintained that is incremented wh
30           * empty list is encountered. This is so more pages are
31           * freed
32           * off fuller lists instead of spinning excessively arou
33           * nd empty
34           * lists
35           */
36          do {

```

```

22     batch_free++;
23     if (++migratetype == MIGRATE_PCPTYPES)
24         migratetype = 0;
25     list = &pcp->lists[migratetype];
26     } while (list_empty(list));
27
28     /* This is the only non-empty list. Free them all. */
29     if (batch_free == MIGRATE_PCPTYPES)
30         batch_free = count;
31
32     do {
33         page = list_last_entry(list, struct page, lru);
34         /* must delete to avoid corrupting pcp list */
35         list_del(&page->lru);
36         pcp->count--;
37
38         if (bulkfree_pcp_prepare(page))
39             continue;
40
41         list_add_tail(&page->lru, &head);
42
43         /*
44         * We are going to put the page back to the glob
45         * pool, prefetch its buddy to speed up later ac
46         * cess
47         * under zone->lock. It is believed the overhead
48         * of
49         * an additional test and calculating buddy_pfn
50         * here
51         * can be offset by reduced memory latency late
52         * r. To
53         * avoid excessive prefetching due to large coun
54         * t, only
55         * prefetch buddy for the first pcp->batch nr of
56         * pages.
57         */
58         if (prefetch_nr++ < pcp->batch)
59             prefetch_buddy(page);
60     } while (--count && --batch_free && !list_empty(list));
61 }
62
63 spin_lock(&zone->lock);
64 isolated_pageblocks = has_isolate_pageblock(zone);
65
66 /*
67 * Use safe version since after __free_one_page(),
68 * page->lru.next will not point to original list.
69 */
70 list_for_each_entry_safe(page, tmp, &head, lru) {
71     int mt = get_pcppage_migratetype(page);
72     /* MIGRATE_ISOLATE page should not go to pcplists */
73     VM_BUG_ON_PAGE(is_migrate_isolate(mt), page);
74     /* Pageblock could have been isolated meanwhile */
75     if (unlikely(isolated_pageblocks))
76         mt = get_pageblock_migratetype(page);
77
78     __free_one_page(page, page_to_pfn(page), zone, 0, mt);
79     trace_mm_page_pcpu_drain(page, 0, mt);
80 }
81 spin_unlock(&zone->lock);
82 }

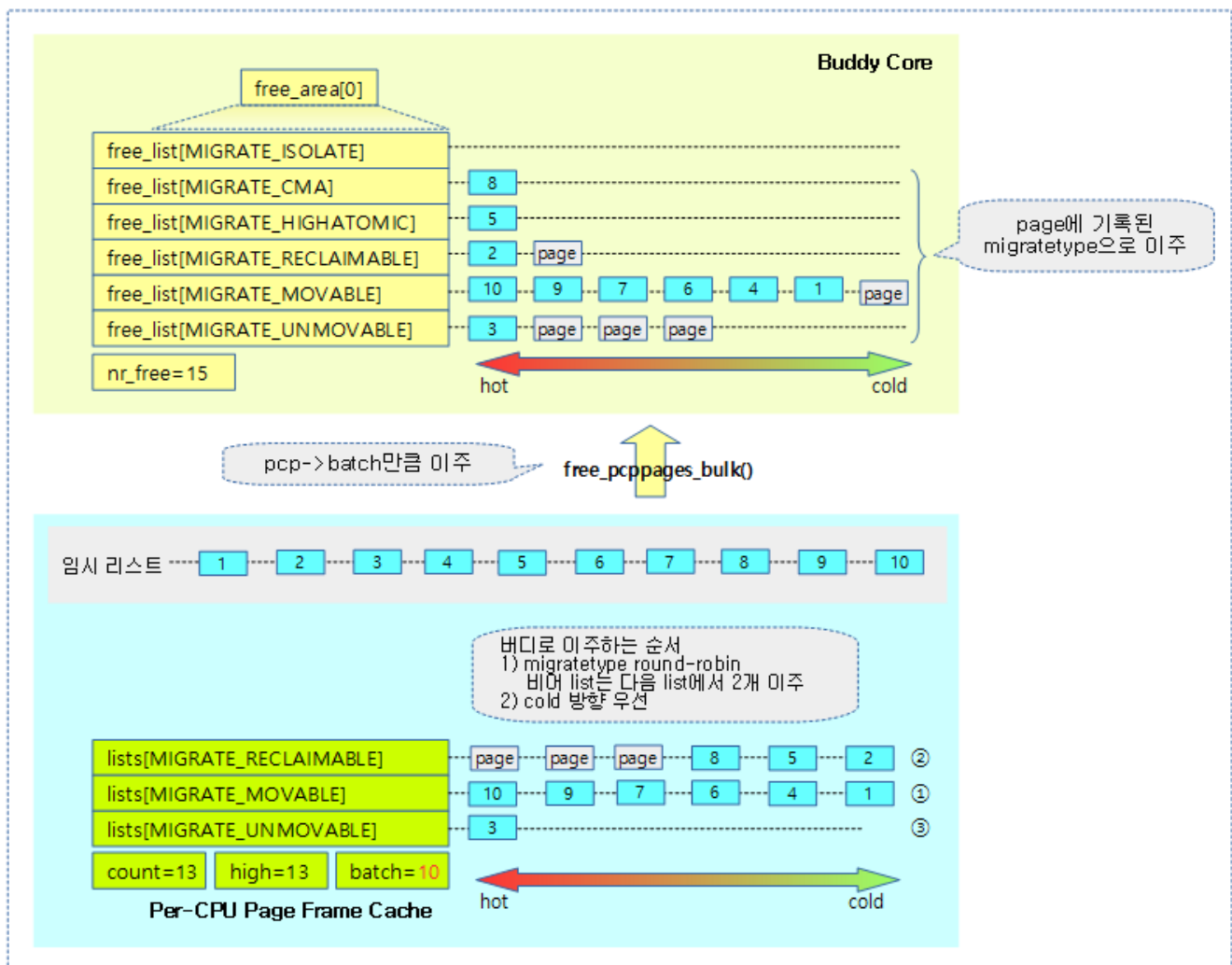
```

The PCP of the request zone is collected by the buddy system as much as it @count.

- Traverse @count number of lines of code 11.

- In lines 21~30 of the code, set the migratetype to traverse the list of pcps of 3 migrate types. However, the empty PCP list should be skipped.
 - At the beginning, proceed with movable(1), reclaimable(2), and unmovable(0) migrate types.
 - Load balance by the number of batch_free, and when the list becomes empty, the batch_free is further increased to prevent too many spins.
 - There are no empty lists, and when you work on three lists, you can rotate one by one.
 - One list is empty, and two are processed when working on the remaining two.
 - If two lists are empty, and only the last list remains, @count is substituted to process them all at once.
- In line 32~41 of code, take the tail direction entry from the pcp list of the specified migratetype and add it to the head direction of the temporary list.
- Lines 52~53 of code from pcp->batch prefetch the buddy page for the page. This is done so that the buddy system can process it with a little faster performance.
- In line 54 of the code, iterate through only batch_free of a list of pcps. However, it is completed when it is empty, or when the @count is 0 and all of them have been processed.
- Traversing the temporary list in code lines 64~74 and using the migrate type to which the page belongs, it retrieves it to the list using the corresponding migrate type of the buddy system. However, if there is an isolate type page in the zone when retrieving, use the migrate type of the page block to which the page belongs.

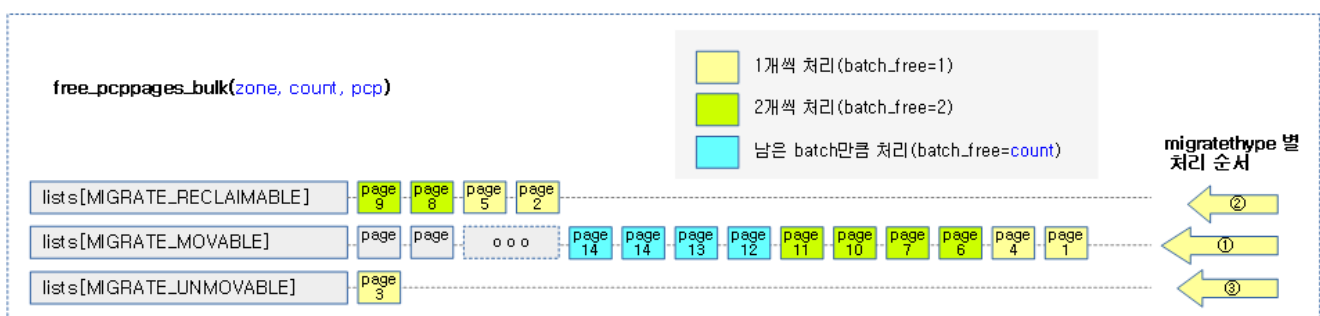
The figure below shows the process and sequence in which PCPs overflow and migrate to buddy by a batch.



(http://jake.dothome.co.kr/wp-content/uploads/2016/04/free_pcppages_bulk-2c.png)

- When a page is migrated to the free_list[0] slot, if a buddy page exists in free_list[0], the buddy page is removed and added by merging it into the next order, free_list[1]. Similarly, if a buddy page is found in free_list[1], it will be merged in the next order, until no buddy page is found.

The following illustration shows the order of pages that are transferred from PCP to Buddy.



(http://jake.dothome.co.kr/wp-content/uploads/2016/04/free_pcppages_bulk-1b.png)

You can check the counter information for PageSets by Zone as follows.

```
01 | pi@pi /proc $ cat zoneinfo
02 | Node 0, zone Normal
03 | pages free 190861
04 | min 2048
```

```

05         low      2560
06         high     3072
07         scanned   0
08         spanned   241664
09         present   241664
10         managed   233403
11         nr_free_pages 190861
12     (...생략...)
13         nr_free_cma 935
14         protection: (0, 0)
15     pagesets
16         cpu: 0
17             count: 50
18             high: 186
19             batch: 31
20     vm stats threshold: 24
21         cpu: 1
22             count: 106
23             high: 186
24             batch: 31
25     vm stats threshold: 24
26         cpu: 2
27             count: 153
28             high: 186
29             batch: 31
30     vm stats threshold: 24
31         cpu: 3
32             count: 156
33             high: 186
34             batch: 31
35     vm stats threshold: 24
36     all_unreclaimable: 0
37     start_pfn: 0
38     inactive_ratio: 1

```

PCP(Per-Cpu Page frame cache) Drain

drain_all_pages()

```

1  /*
2   * Spill all the per-cpu pages from all CPUs back into the buddy allocat
3   or.
4   *
5   * When zone parameter is non-NULL, spill just the single zone's pages.
6   *
7   * Note that this can be extremely slow as the draining happens in a wor
8   kqueue.
9   */
10
11 void drain_all_pages(struct zone *zone)
12 {
13     int cpu;
14
15     /*
16      * Allocate in the BSS so we wont require allocation in
17      * direct reclaim path for CONFIG_CPUMASK_OFFSTACK=y
18      */
19     static cpumask_t cpus_with_pcps;
20
21     /*
22      * Make sure nobody triggers this path before mm_percpu_wq is fu
23      lly
24      * initialized.
25      */

```

```

15     if (WARN_ON_ONCE(!mm_percpu_wq))
16         return;
17
18     /*
19      * Do not drain if one is already in progress unless it's specif
20      * a zone. Such callers are primarily CMA and memory hotplug and
21      * the drain to be complete when the call returns.
22      */
23     if (unlikely(!mutex_trylock(&pcpu_drain_mutex))) {
24         if (!zone)
25             return;
26         mutex_lock(&pcpu_drain_mutex);
27     }
28
29     /*
30      * We don't care about racing with CPU hotplug event
31      * as offline notification will cause the notified
32      * cpu to drain that CPU pcps and on_each_cpu_mask
33      * disables preemption as part of its processing
34      */
35     for_each_online_cpu(cpu) {
36         struct per_cpu_pageset *pcp;
37         struct zone *z;
38         bool has_pcps = false;
39
40         if (zone) {
41             pcp = per_cpu_ptr(zone->pageset, cpu);
42             if (pcp->pcp.count)
43                 has_pcps = true;
44         } else {
45             for_each_populated_zone(z) {
46                 pcp = per_cpu_ptr(z->pageset, cpu);
47                 if (pcp->pcp.count) {
48                     has_pcps = true;
49                     break;
50                 }
51             }
52         }
53
54         if (has_pcps)
55             cpumask_set_cpu(cpu, &cpus_with_pcps);
56         else
57             cpumask_clear_cpu(cpu, &cpus_with_pcps);
58     }
59
60     for_each_cpu(cpu, &cpus_with_pcps) {
61         struct pcpu_drain *drain = per_cpu_ptr(&pcpu_drain, cp
62         u);
63
64         drain->zone = zone;
65         INIT_WORK(&drain->work, drain_local_pages_wq);
66         queue_work_on(cpu, mm_percpu_wq, &drain->work);
67     }
68     for_each_cpu(cpu, &cpus_with_pcps)
69         flush_work(&per_cpu_ptr(&pcpu_drain, cpu)->work);
70     mutex_unlock(&pcpu_drain_mutex);
71 }

```

Move the Per-CPU Page Frame Cache on all online CPUs in the specified zone to the buddy memory allocator. If no zone is specified, this is done for all populated zones.

drain_local_pages()

mm/page_alloc.c

```

1  /*
2  * Spill all of this CPU's per-cpu pages back into the buddy allocator.
3  *
4  * The CPU has to be pinned. When zone parameter is non-NULL, spill just
5  * the single zone's pages.
6  */

1  void drain_local_pages(struct zone *zone)
2  {
3      int cpu = smp_processor_id();
4
5      if (zone)
6          drain_pages_zone(cpu, zone);
7      else
8          drain_pages(cpu);
9  }

```

drain_pages()

mm/page_alloc.c

```

1  /*
2  * Drain pcplists of all zones on the indicated processor.
3  *
4  * The processor must either be the current processor and the
5  * thread pinned to the current processor or a processor that
6  * is not online.
7  */

1  static void drain_pages(unsigned int cpu)
2  {
3      struct zone *zone;
4
5      for_each_populated_zone(zone) {
6          drain_pages_zone(cpu, zone);
7      }
8  }

```

Empty the Per-CPU Page Fram Cache for all enabled zones.

drain_pages_zone()

mm/page_alloc.c

```

1  /*
2  * Drain pcplists of the indicated processor and zone.
3  *
4  * The processor must either be the current processor and the
5  * thread pinned to the current processor or a processor that
6  * is not online.
7  */

01 static void drain_pages_zone(unsigned int cpu, struct zone *zone)
02 {
03     unsigned long flags;
04     struct per_cpu_pageset *pset;
05     struct per_cpu_pages *pcp;
06
07     local_irq_save(flags);
08     pset = per_cpu_ptr(zone->pageset, cpu);

```

```

09
10     pcp = &pset->pcp;
11     if (pcp->count)
12         free_pcppages_bulk(zone, pcp->count, pcp);
13     local_irq_restore(flags);
14 }

```

All pages registered in the Per-CPU Page Fram Cache for the request zone are migrated to the buddy system.

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 - 현재 글
- Zoned Allocator -6- (Watermark) (<http://jake.dothome.co.kr/zoned-allocator-watermark>) | 문c
- 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
- setup_per_cpu_pageset() ([http://jake.dothome.co.kr/setup_per_cpu_pageset\(\)](http://jake.dothome.co.kr/setup_per_cpu_pageset())) | Qc
- page_alloc_init() (http://jake.dothome.co.kr/page_alloc_init) | Qc
- CPU Bitmap (API) (<http://jake.dothome.co.kr/cpu-api>) | Qc

4 thoughts to "Zoned Allocator -5- (Per-CPU Page Frame

Cache)"

**TARO**2021-05-20 22:06 (<http://jake.dothome.co.kr/per-cpu-page-frame-cache/#comment-305337>)

Hello

In

the `__rmqueue_pcplist` function, `page = list_first_entry(list, struct page, lru)` to get page.

If you look at the `list_first_entry` function, you get the page from the first entry in the list using the offset of LRU,

which is the concatenation list inside the struct page.

If the list is created by the `setup_per_cpu_pageset()` function, is this list the same as the LRU of the struct page?

RESPONSE (/PER-CPU-PAGE-FRAME-CACHE/?REPLYTOCOM=305337#RESPOND)

**MOON YOUNG-IL (HTTP://JAKE.DOTHOME.CO.KR)**2021-05-22 14:16 (<http://jake.dothome.co.kr/per-cpu-page-frame-cache/#comment-305338>)

Hello? Mr. Taro.

The last question sentence "Is this list the same as the LRU on the struct page?" Except for that,

you solved it correctly. The last sentence of the question makes it difficult for me to grasp the exact intent of the question.

The list is the `lists[]` of the `per_cpu_pages` struct member, and the pages are hung from this list.

```
list = &pcp->lists[migratetype];
```

If you haven't figured it out yet, I'd love to hear you clarify the last question a little more closely.

RESPONSE (/PER-CPU-PAGE-FRAME-CACHE/?REPLYTOCOM=305338#RESPOND)

**TARO**2021-05-23 16:34 (<http://jake.dothome.co.kr/per-cpu-page-frame-cache/#comment-305342>)

What

I don't really understand is the part in your answer that says 'pages are stuck in `lists[]`'.

When a page hangs on lists[], is the address of the struct list_head lru inside the page hanging?

page = list_first_entry(list, struct page, lru) 를 통해 page를 가져올 때
list_first_entry(list, struct page, lru)
list_entry(list->next, struct page, lru)
container_of(list->next, struct page, lru) 이고,

In the end, it becomes (struct page *) (list->next - offsetof(struct page, lru)) and we get the address of the page by subtracting the offset of the LRU inside the struct page.

그러려면, lists[]에 저장되는 것이 struct page 내부의 lru의 주소여야 한다고 생각하는데 제가 생각하는 것이 맞을까요?

응답 (/PER-CPU-PAGE-FRAME-CACHE/?REPLYTOCOM=305342#RESPOND)



문영일 (HTTP://JAKE.DOTHOME.CO.KR)

2021-05-24 14:11 (<http://jake.dothome.co.kr/per-cpu-page-frame-cache/#comment-305349>)

네. 맞습니다.

list에 두 개의 페이지가 연결되어 있는 경우 다음과 같이 연결됩니다.

(참고로 lru는 double linked list로, list_head 구조체이며 *next와 *prev 멤버가 사용됩니다)

list < ----> A page->lru < ----> B page->lru

위를 더 자세히 표현하면 다음과 같습니다.

list->next —> A page->lru.next —> B page->lru.next —> 맨 앞 list->next

list->prev < --- A page->lru.prev < --- B page->lru.prev < --- 맨 앞 list->prev

감사합니다.

응답 (/PER-CPU-PAGE-FRAME-CACHE/?REPLYTOCOM=305349#RESPOND)

댓글 남기기

이메일은 공개되지 않습니다. 필수 입력창은 * 로 표시되어 있습니다

댓글

이름 *

이메일 *

웹사이트

[댓글 작성](#)[◀ Control Group for Memory \(http://jake.dothome.co.kr/mem_cgroup/\)](http://jake.dothome.co.kr/mem_cgroup/)[Zoned Allocator -10- \(LRU & pagevecs\) ▶ \(http://jake.dothome.co.kr/lru-lists-pagevecs/\)](http://jake.dothome.co.kr/lru-lists-pagevecs/)[문c 블로그 \(2015 ~ 2023\)](#)