# Zoned Allocator -8- (Direct Compact-Isolation)

📅 2016-07-05 (http://jake.dothome.co.kr/zonned-allocator-isolation/)    👤 Moon Young-il (http://jake.dothome.co.kr/author/admin/)    📁 Linux Kernel (http://jake.dothome.co.kr/category/linux/)
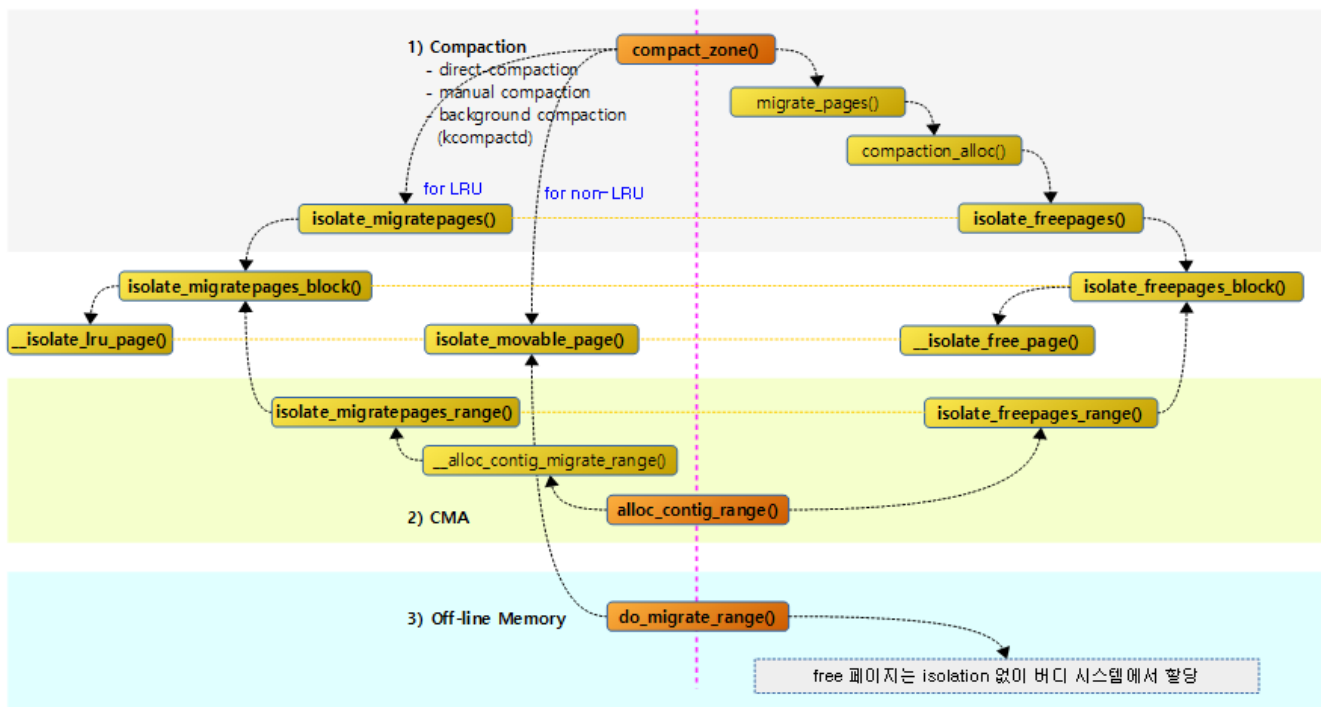
<div align="right">&lt;kernel v5.0&gt;</div>

## Zoned Allocator -8- (Direct Compact-Isolation)

Isolation can be used for the following purposes:

- Compaction
    - Migrate migrating migratable pages to the upper part of the zone to make up for the missing high order pages.
    - Direct-compaction, manual-compaction, kcompactd
- Off-line Memory
    - Migrate all in-use movable pages located in the memory area to be offline to another area.
- CMA
    - In order to secure a continuous physical space within the scope requested by the CMA area, the movable pages that are temporarily occupying (?) in the CMA area are migrated.

The following diagram shows the main function call process for isolation in two parts. On the left is the process of isolating the page to be moved, and on the right is the process of isolating or securing the free page.
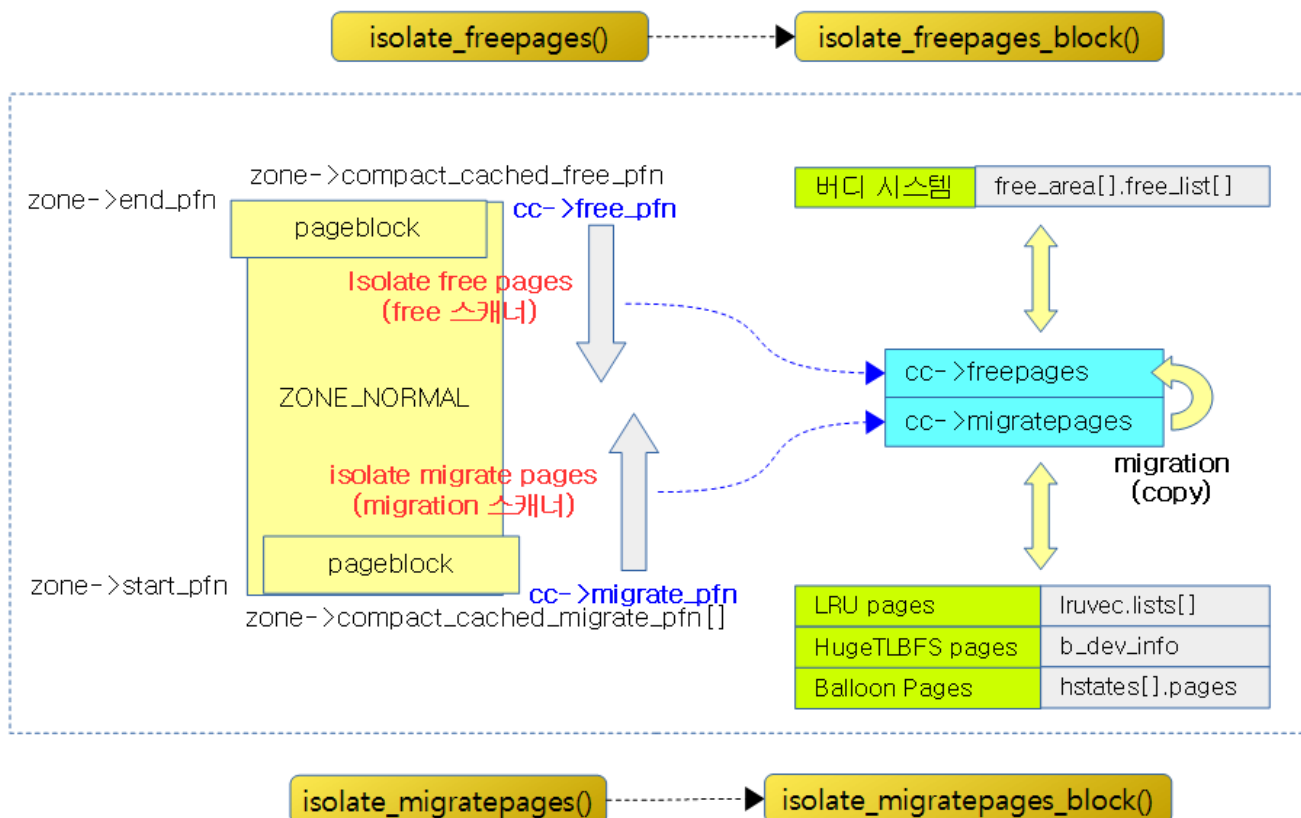


(http://jake.dothome.co.kr/wp-content/uploads/2016/07/compaction-6b.png)

## Two scanners for Compaction

The following diagram shows the direction in which the two migration scanners and the free scanner work when performing a compaction, isolating the pages to be used for their respective purposes.

- The migration scanner scans page blocks upwards to get migratable pages.
  - isolate_migratepages_block() isolates the migratable page from a single selected page block and adds it to the cc->migratepages list.
- The free scanner scans the page blocks in a downward direction to get the free pages.
  - isolate_freepages_block() isolates a free page of type request from a single selected page block and adds it to the cc->freepages list.



(http://jake.dothome.co.kr/wp-content/uploads/2016/07/isolation-1.png)

---

# Migrate Scanner

## Isolate migratepages

## isolate_migratepages()

mm/compaction.c

```
1  /*
2   * Isolate all pages that can be migrated from the first suitable block,
3   * starting at the block pointed to by the migrate scanner pfn within
4   * compact_control.
5   */
```

```
01   static isolate_migrate_t isolate_migratepages(struct zone *zone,
02                                           struct compact_control *cc)
03   {
04           unsigned long block_start_pfn;
05           unsigned long block_end_pfn;
06           unsigned long low_pfn;
07           struct page *page;
08           const isolate_mode_t isolate_mode =
09                   (sysctl_compact_unevictable_allowed ? ISOLATE_UNEVICTABL
E : 0) |
10                   (cc->mode != MIGRATE_SYNC ? ISOLATE_ASYNC_MIGRATE : 0);
11
12           /*
13            * Start at where we last stopped, or beginning of the zone as
14            * initialized by compact_zone()
15            */
16           low_pfn = cc->migrate_pfn;
17           block_start_pfn = pageblock_start_pfn(low_pfn);
18           if (block_start_pfn < zone->zone_start_pfn)
19                   block_start_pfn = zone->zone_start_pfn;
20
21           /* Only scan within a pageblock boundary */
22           block_end_pfn = pageblock_end_pfn(low_pfn);
23
24           /*
25            * Iterate over whole pageblocks until we find the first suitabl
e.
26            * Do not cross the free scanner.
27            */
28           for (; block_end_pfn <= cc->free_pfn;
29                           low_pfn = block_end_pfn,
30                           block_start_pfn = block_end_pfn,
31                           block_end_pfn += pageblock_nr_pages) {
32
33                   /*
34                    * This can potentially iterate a massively long zone wi
th
35                    * many pageblocks unsuitable, so periodically check if
we
36                    * need to schedule, or even abort async compaction.
37                    */
38                   if (!(low_pfn % (SWAP_CLUSTER_MAX * pageblock_nr_pages))
39                                           && compact_should_abort
(cc))
40                           break;
41
42                   page = pageblock_pfn_to_page(block_start_pfn, block_end_
pfn,
43                                                                         z
one);
44                   if (!page)
45                           continue;
46
47                   /* If isolation recently failed, do not retry */
48                   if (!isolation_suitable(cc, page))
49                           continue;
50
51                   /*
52                    * For async compaction, also only scan in MOVABLE block
s.
53                    * Async compaction is optimistic to see if the minimum
amount
54                    * of work satisfies the allocation.
55                    */
56                   if (!suitable_migration_source(cc, page))
57                           continue;
58
59                   /* Perform the isolation */
```

```
60              low_pfn = isolate_migratepages_block(cc, low_pfn,
61                                         block_end_pfn, isolate_m
   ode);
62
63              if (!low_pfn || cc->contended)
64                  return ISOLATE_ABORT;
65
66              /*
67               * Either we isolated something and proceed with migrati
   on. Or
68               * we failed and compact_zone should decide if we should
69               * continue or not.
70               */
71              break;
72          }
73
74          /* Record where migration scanner will be restarted. */
75          cc->migrate_pfn = low_pfn;
76
77          return cc->nr_migratepages ? ISOLATE_SUCCESS : ISOLATE_NONE;
78      }
```
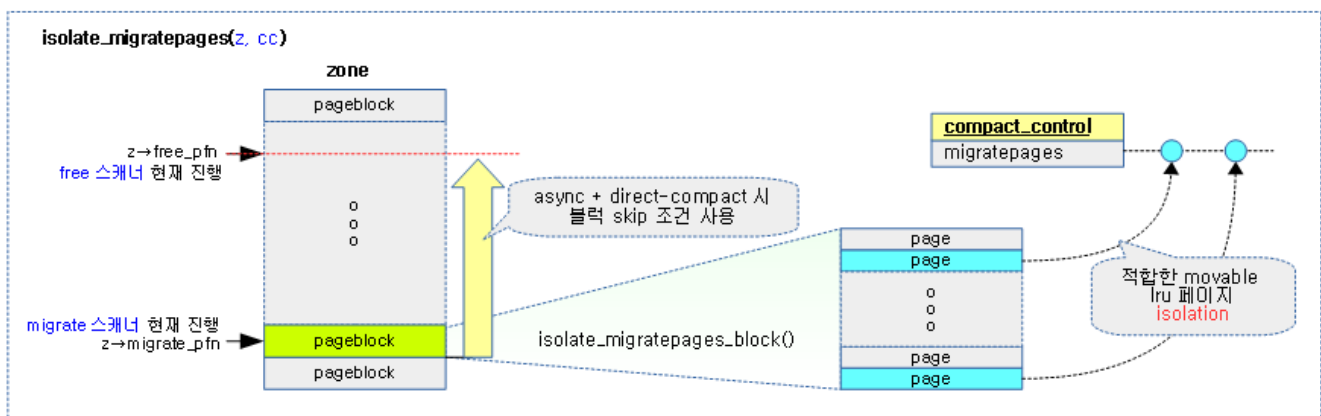
Isolate migratable pages in the requested zone.

- Determine the isolate mode from code lines 8~10.
    - ISOLATE_UNMAPPED(0x2)
        - Only unmapped pages support isolation, and mapped files (file cache, swap cache, etc.) are not processed.
    - ISOLATE_ASYNC_MIGRATE(0x4)
        - Briefly perform isolation and migration. This is done only for page blocks with a migrate type of more than 50% of the blocks, and if 1 order isolation is possible, it will stop without proceeding any further. In addition, it also stops preempt requests.
    - ISOLATE_UNEVICTABLE(0x8)
        - Even unevictable pages have isolation performed.
- In code lines 16~22, find the pfn of the beginning and end of the block to be scanned with the pfn pointing to the location where the migrate scanner wants to scan upwards.
    - e.g. migrate_pfn=0x250, pageblock_order=9
        - block_start_pfn=0x200, block_end_pfn=0x400
- Traverse the Migrate scanner block by block from code lines 28~31 to the location of the free scanner that scans downward.
- Since there is a large scope to be processed in line 38~40 of the code, periodically check for interruptions.
    - Check Frequency
        - SWAP_CLUSTER_MAX (32) page blocks
    - Interrupting Elements
        - Async migration is being processed, but there is a preemption request from a higher priority task.
- In lines 42~45 of the code, get the first page of the page block in the scope of the request zone. If the page leaves the request zone, it returns null and skips.
- In lines 48~49 of code, if you don't want to isolate the block in that zone, skip it.
    - When direct-compact is run in async mode, it returns false to skip the block if the skip bit that is set if isolation has recently failed on that pageblock.

- In line 56~57 of code, when operating in asynchronous direct-compact mode, if the block migrate type is not the same as the requested migrate type, skip it. However, if it is a synchronous (migrate_sync*) request, or if it is a manual and kcompactd request, it will always return true to make an unconditional attempt to isolate the block.
- In code lines 60~61, only the migratable pages are isolated from the pages in that block.
- If isolation fails in line 63~64 or has to be aborted halfway through, it will exit with ISOLATE_ABORT result.
- If isolation completes normally on line 71 of code, it exits the loop.
- Remembers the current position of the migrate scanner in process from code lines 75~77 and returns an isolation result.

The following figure shows how the Migrate scanner uses the isolate_migratepages() function to isolate movable LRU pages against page blocks with conformity conditions.

- lru page isolation (X) -> migratable page suitable movable



(http://jake.dothome.co.kr/wp-content/uploads/2016/07/isolate_migratepages-1.png)

# Isolate migratepages block

The following functions are called and used for two purposes:

- Compaction
  - The isolate_migratepages() function can be called with the addition of the following isolate mode-related flags:
    - If you are not in sync compaction mode, a ISOLATE_ASYNC_MIGRATE is added.
    - If sysctl_compact_unevictable_allowed (default=1) is set, a ISOLATE_UNEVICTABLE mode is added.
- CMA
  - isolate_migratepages_range() function is called in ISOLATE_UNEVICTABLE mode.

### isolate_migratepages_block()

mm/compaction.c -1/4-

```
01  /**
02   * isolate_migratepages_block() - isolate all migrate-able pages within
03   *                                 a single pageblock
04   * @cc:         Compaction control structure.
05   * @low_pfn:    The first PFN to isolate
06   * @end_pfn:    The one-past-the-last PFN to isolate, within same pagebl
ock
07   * @isolate_mode: Isolation mode to be used.
08   *
09   * Isolate all pages that can be migrated from the range specified by
10   * [low_pfn, end_pfn). The range is expected to be within same pagebloc
k.
11   * Returns zero if there is a fatal signal pending, otherwise PFN of the
12   * first page that was not scanned (which may be both less, equal to or
more
13   * than end_pfn).
14   *
15   * The pages are isolated on cc->migratepages list (not required to be e
mpty),
16   * and cc->nr_migratepages is updated accordingly. The cc->migrate_pfn f
ield
17   * is neither read nor updated.
18   */
```

```
01  isolate_migratepages_block(struct compact_control *cc, unsigned long low
_pfn,
02                          unsigned long end_pfn, isolate_mode_t isolate_mo
de)
03  {
04          struct zone *zone = cc->zone;
05          unsigned long nr_scanned = 0, nr_isolated = 0;
06          struct lruvec *lruvec;
07          unsigned long flags = 0;
08          bool locked = false;
09          struct page *page = NULL, *valid_page = NULL;
10          unsigned long start_pfn = low_pfn;
11          bool skip_on_failure = false;
12          unsigned long next_skip_pfn = 0;
13
14          /*
15           * Ensure that there are not too many pages isolated from the LR
U
16           * list by either parallel reclaimers or compaction. If there ar
e,
17           * delay for some time until fewer pages are isolated
18           */
19          while (unlikely(too_many_isolated(zone))) {
20                  /* async migration should just abort */
21                  if (cc->mode == MIGRATE_ASYNC)
22                          return 0;
23
24                  congestion_wait(BLK_RW_ASYNC, HZ/10);
25
26                  if (fatal_signal_pending(current))
27                          return 0;
28          }
29
30          if (compact_should_abort(cc))
31                  return 0;
32
33          if (cc->direct_compaction && (cc->mode == MIGRATE_ASYNC)) {
34                  skip_on_failure = true;
35                  next_skip_pfn = block_end_pfn(low_pfn, cc->order);
36          }
37
38          /* Time to isolate some pages for migration */
39          for (; low_pfn < end_pfn; low_pfn++) {
40
```

```
41          if (skip_on_failure && low_pfn >= next_skip_pfn) {
42                  /*
43                   * We have isolated all migration candidates in the
44                   * previous order-aligned block, and did not skip it due
45                   * to failure. We should migrate the pages now and
46                   * hopefully succeed compaction.
47                   */
48                  if (nr_isolated)
49                      break;
50
51                  /*
52                   * We failed to isolate in the previous order-aligned
53                   * block. Set the new boundary to the end of the
54                   * current block. Note we can't simply increase
55                   * next_skip_pfn by 1 << order, as low_pfn might have
56                   * been incremented by a higher number due to skipping
57                   * a compound or a high-order buddy page in the
58                   * previous loop iteration.
59                   */
60                  next_skip_pfn = block_end_pfn(low_pfn, cc->order);
61          }
```

Isolate migratable pages in a one-page block of the requested migratetype.

- In line 5 of the code, initialize the number of scanned pages and the number of isolated pages to zero.
- In line 19~28 of the code, there is a small probability that only if there are too many isolated pages, it will loop around and delay processing by 100ms. However, in the following situations, it aborts processing and returns 0.
    - During asynchronous migration, excessive processing is abandoned.
    - If the processing of the SIGKILL signal is delayed, the function will be aborted.
    - 참고: Avoid the use of congestion_wait under zone pressure (https://lwn.net/Articles/377709/)
- In line 30~31 of the code, if there is a preemption request from a higher priority task during asynchronous processing, the function processing is aborted.
- In line 33~36 of code, if you are processing async migration using direct-compaction, set the skip_on_failure to true.
    - skip_on_failure
        - In order to quickly assign the requested order, the page is searched by order rather than a one-page block, and if there is any isolation, it is terminated.
- Travers the pages in the range to be processed one by one in line 39~61 of the code, but if the skip_on_failure is set and the processing by order unit is completed, specify the end pfn of the next order. If there is any isolated page in the block, it exits the loop.

mm/compaction.c -2/4-

```
01  .               /*
```

```
02                              * Periodically drop the lock (if held) regardless of it
                        s
03                              * contention, to give chance to IRQs. Abort async compa
                        ction
04                              * if contended.
05                              */
06                             if (!(low_pfn % SWAP_CLUSTER_MAX)
07                                 && compact_unlock_should_abort(zone_lru_lock(zone),
                        flags,
08                                                                          &locked,
                        cc))
09                                     break;
10
11                             if (!pfn_valid_within(low_pfn))
12                                     goto isolate_fail;
13                             nr_scanned++;
14
15                             page = pfn_to_page(low_pfn);
16
17                             if (!valid_page)
18                                     valid_page = page;
19
20                             /*
21                              * Skip if free. We read page order here without zone lo
                        ck
22                              * which is generally unsafe, but the race window is sma
                        ll and
23                              * the worst thing that can happen is that we skip some
24                              * potential isolation targets.
25                              */
26                             if (PageBuddy(page)) {
27                                     unsigned long freepage_order = page_order_unsafe
                        (page);
28
29                                     /*
30                                      * Without lock, we cannot be sure that what we
                        got is
31                                      * a valid page order. Consider only values in t
                        he
32                                      * valid order range to prevent low_pfn overflo
                        w.
33                                      */
34                                     if (freepage_order > 0 && freepage_order < MAX_O
                        RDER)
35                                             low_pfn += (1UL << freepage_order) - 1;
36                                     continue;
37                             }
38
39                             /*
40                              * Regardless of being on LRU, compound pages such as TH
                        P and
41                              * hugetlbfs are not to be compacted. We can potentially
                        save
42                              * a lot of iterations if we skip them at once. The chec
                        k is
43                              * racy, but we can consider only valid values and the o
                        nly
44                              * danger is skipping too much.
45                              */
46                             if (PageCompound(page)) {
47                                     const unsigned int order = compound_order(page);
48
49                                     if (likely(order < MAX_ORDER))
50                                             low_pfn += (1UL << order) - 1;
51                                     goto isolate_fail;
52                             }
53
54                             /*
```

```
55              * Check may be lockless but that's ok as we recheck lat
   er.
56              * It's possible to migrate LRU and non-lru movable page
   s.
57              * Skip any other type of page
58              */
59             if (!PageLRU(page)) {
60                     /*
61                      * __PageMovable can return false positive so we
   need
62                      * to verify it under page_lock.
63                      */
64                     if (unlikely(__PageMovable(page)) &&
65                                 !PageIsolated(page)) {
66                             if (locked) {
67                                     spin_unlock_irqrestore(zone_lru_
   lock(zone),
68                                                                      f
   lags);
69                                     locked = false;
70                             }
71
72                             if (!isolate_movable_page(page, isolate_
   mode))
73                                     goto isolate_success;
74                     }
75
76                     goto isolate_fail;
77             }
```

- Since there is a large scope to be processed in line 6~9 of the code, we try to increase the irq latency by periodically unlocking it in this routine. Also, if there is an interruption element, it aborts processing and exits the loop.
    - Check Frequency
        - SWAP_CLUSTER_MAX (32) pages per block
    - Interrupting Elements
        - The SIGKILL signal is delayed.
        - If there is a preemption request from a high-priority task during an asynchronous migration.
- Update the first valid page from code lines 11~18, and increase the scan counter. If the page is not valid, go to the isolate_fail label.
- In code lines 26~37, the migrate scanner migrates only movable pages that are in use. Therefore, if it is a free page managed by the buddy system, it will be skipped.
- In lines 46~52 of code, the compound(slab, hugetlbfs, thp) page also has no compaction effect, so it goes to the isolate_fail label.
- If the LRU page is not user-assigned in code lines 59~77, go to the isolate_fail label. However, if it is a non-LRU movable page that has not already been isolated, it will be isolated, and if it is successful, it will be moved to the isolate_success label.
    - GPU, ZSRAM (z3fold, zsmalloc), and balloon drivers support migration of non-LRU-movable pages.

mm/compaction.c -3/4-

```
01              /*
```

```
02                          * Migration will fail if an anonymous page is pinned in
    memory,
03                          * so avoid taking lru_lock and isolating it unnecessari
    ly in an
04                          * admittedly racy check.
05                          */
06                         if (!page_mapping(page) &&
07                             page_count(page) > page_mapcount(page))
08                                 goto isolate_fail;
09
10                         /*
11                          * Only allow to migrate anonymous pages in GFP_NOFS con
    text
12                          * because those do not depend on fs locks.
13                          */
14                         if (!(cc->gfp_mask & __GFP_FS) && page_mapping(page))
15                                 goto isolate_fail;
16
17                         /* If we already hold the lock, we can skip some recheck
    ing */
18                         if (!locked) {
19                                 locked = compact_trylock_irqsave(zone_lru_lock(z
    one),
20                                                                     &flags,
    cc);
21                                 if (!locked)
22                                         break;
23
24                                 /* Recheck PageLRU and PageCompound under lock
    */
25                                 if (!PageLRU(page))
26                                         goto isolate_fail;
27
28                                 /*
29                                  * Page become compound since the non-locked che
    ck,
30                                  * and it's on LRU. It can only be a THP so the
    order
31                                  * is safe to read and it's 0 for tail pages.
32                                  */
33                                 if (unlikely(PageCompound(page))) {
34                                         low_pfn += (1UL << compound_order(page))
    - 1;
35                                         goto isolate_fail;
36                                 }
37                         }
38
39                         lruvec = mem_cgroup_page_lruvec(page, zone->zone_pgdat);
40
41                         /* Try isolate the page */
42                         if (__isolate_lru_page(page, isolate_mode) != 0)
43                                 goto isolate_fail;
44
45                         VM_BUG_ON_PAGE(PageCompound(page), page);
46
47                         /* Successfully isolated */
48                         del_page_from_lru_list(page, lruvec, page_lru(page));
49                         inc_node_page_state(page,
50                                         NR_ISOLATED_ANON + page_is_file_cache(pa
    ge));
51
52  isolate_success:
53                         list_add(&page->lru, &cc->migratepages);
54                         cc->nr_migratepages++;
55                         nr_isolated++;
56
57                         /*
```

```
58                          * Record where we could have freed pages by migration a
        nd not
59                          * yet flushed them to buddy allocator.
60                          * - this is the lowest page that was isolated and likel
        y be
61                          * then freed by migration.
62                          */
63                         if (!cc->last_migrated_pfn)
64                                 cc->last_migrated_pfn = low_pfn;
65
66                         /* Avoid isolating too much */
67                         if (cc->nr_migratepages == COMPACT_CLUSTER_MAX) {
68                                 ++low_pfn;
69                                 break;
70                         }
71
72                         continue;
```

- In line 6~8 of code, move the anonymous page to the isolate_fail if it is in use in a different kernel context.
- If lines 14~15 of the code are mapped to the file system, but fs is not available using GFP_NOFS, go to the isolate_fail label.
- In code lines 18~37, the lock was released from above at regular intervals. In this case, you will need to reacquire it and check the page conditions again. If it's not an LRU page, or if it's a Compound page, go to the isolate_fail label.
- In line 39 of code, select either the node's LRU or the node's LRU from MEMCG for the LRU list.
- In lines 42~50 of code, detach the page from the LRU list and increment the NR_ISOLATED_ANON or NR_ISOLATE_FILE counters.
- In code lines 52~55, isolation is considered successful and can be moved: isolate_success: label. Add the separated pages to the migratepages list and increment the associated stats.
- Update code lines 63~64 if the last migrated page is never specified.
- In line 67~70 of the code, if the migrate page reaches the appropriate amount (32), the processing is stopped.
- Continue on line 72 of code and loop around to process the next page.


mm/compaction.c -4/4-

```
01  isolate_fail:
02                         if (!skip_on_failure)
03                                 continue;
04
05                         /*
06                          * We have isolated some pages, but then failed. Release
        them
07                          * instead of migrating, as we cannot form the cc->order
        buddy
08                          * page anyway.
09                          */
10                         if (nr_isolated) {
11                                 if (locked) {
12                                         spin_unlock_irqrestore(zone_lru_lock(zon
        e), flags);
13                                         locked = false;
14                                 }
15                                 putback_movable_pages(&cc->migratepages);
16                                 cc->nr_migratepages = 0;
17                                 cc->last_migrated_pfn = 0;
```

```
18                                  nr_isolated = 0;
19                          }
20
21                          if (low_pfn < next_skip_pfn) {
22                                  low_pfn = next_skip_pfn - 1;
23                                  /*
24                                   * The check near the loop beginning would have
     updated
25                                   * next_skip_pfn too, but this is a bit simpler.
26                                   */
27                                  next_skip_pfn += 1UL << cc->order;
28                          }
29                  }
30
31                  /*
32                   * The PageBuddy() check could have potentially brought us outsi
     de
33                   * the range to be scanned.
34                   */
35                  if (unlikely(low_pfn > end_pfn))
36                          low_pfn = end_pfn;
37
38                  if (locked)
39                          spin_unlock_irqrestore(zone_lru_lock(zone), flags);
40
41                  /*
42                   * Update the pageblock-skip information and cached scanner pfn,
43                   * if the whole pageblock was scanned without isolating any pag
     e.
44                   */
45                  if (low_pfn == end_pfn)
46                          update_pageblock_skip(cc, valid_page, nr_isolated, tru
     e);
47
48                  trace_mm_compaction_isolate_migratepages(start_pfn, low_pfn,
49                                                  nr_scanned, nr_isolate
     d);
50
51                  cc->total_migrate_scanned += nr_scanned;
52                  if (nr_isolated)
53                          count_compact_events(COMPACTISOLATED, nr_isolated);
54
55                  return low_pfn;
56  }
```
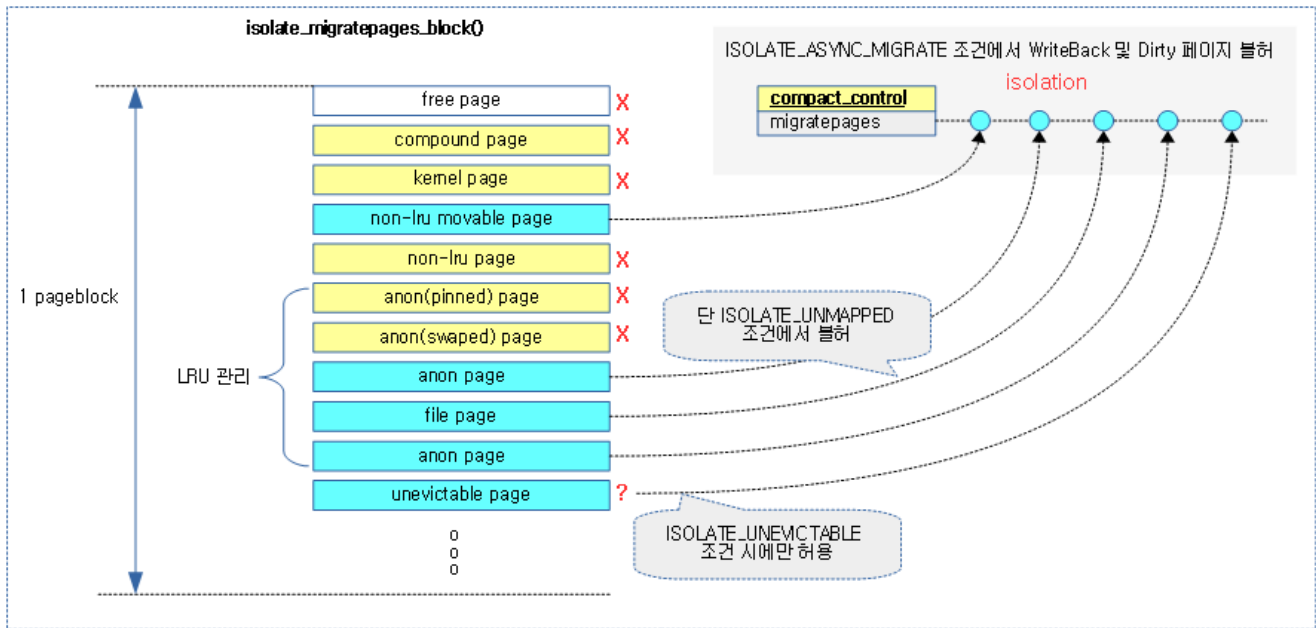
- In code lines 1~3, isolate_faile: is a label. If isolation fails, if the skip_on_failure is false, continue processing the next page.
- Revert the pages that were isolated in code lines 10~19 and invalidate them.
  - If skip_on_failure is true, if the order unit pages fail in isolation, try the next order unit page.
- In lines 21~28 of the code, change the low_pfn to the next order unit page.
- In lines 35~36 of the code, the low_pfn is restricted from going beyond the end of the processing range.
- If the lock is applied in line 38~39 of the code, unlock it.
- If the end of the page block is processed from line 45~46 of the code and there are no isolated pages processed, set the migrate skip bit of the page block corresponding to valid_page to 1 so that it will be skipped in the next scan. Then, set the start pfn of the migrate scanner (both async and sync) to that page.
- Update the scan counter and isolation counter on lines 51~53 of code.

The following image shows the isolate_migratepages_block() function to isolate a migratable page.



(http://jake.dothome.co.kr/wp-content/uploads/2016/07/isolate_migratepages-block-1b.png)

## too_many_isolated()

mm/compaction.c

```
01  /* Similar to reclaim, but different enough that they don't share logic
    */
02  static bool too_many_isolated(struct zone *zone)
03  {
04          unsigned long active, inactive, isolated;
05
06          inactive = zone_page_state(zone, NR_INACTIVE_FILE) +
07                                      zone_page_state(zone, NR_INACTIV
    E_ANON);
08          active = zone_page_state(zone, NR_ACTIVE_FILE) +
09                                      zone_page_state(zone, NR_ACTIVE_
    ANON);
10          isolated = zone_page_state(zone, NR_ISOLATED_FILE) +
11                                      zone_page_state(zone, NR_ISOLATE
    D_ANON);
12
13          return isolated > (inactive + active) / 2;
14  }
```

isolated returns true if the number of pages is too high.

- file(active+inactive) true if more than half of the pages

## compact_should_abort()

mm/compaction.c

```
1  /*
2   * Aside from avoiding lock contention, compaction also periodically che
   cks
```

```
 3   * need_resched() and either schedules in sync compaction or aborts async
     c
 4   * compaction. This is similar to what compact_unlock_should_abort() doe
     s, but
 5   * is used where no lock is concerned.
 6   *
 7   * Returns false when no scheduling was needed, or sync compaction sched
     uled.
 8   * Returns true when async compaction should abort.
 9   */
```

```
01  static inline bool compact_should_abort(struct compact_control *cc)
02  {
03          /* async compaction aborts if contended */
04          if (need_resched()) {
05                  if (cc->mode == MIGRATE_ASYNC) {
06                          cc->contended = true;
07                          return true;
08                  }
09
10                  cond_resched();
11          }
12
13          return false;
14  }
```

If there is a lisse-scheduled request from a high-priority task while the asynchronous migration is being processed, true is returned.

## compact_unlock_should_abort()

mm/compaction.c

```
01  /*
02   * Compaction requires the taking of some coarse locks that are potentia
     lly
03   * very heavily contended. The lock should be periodically unlocked to a
     void
04   * having disabled IRQs for a long time, even when there is nobody waiti
     ng on
05   * the lock. It might also be that allowing the IRQs will result in
06   * need_resched() becoming true. If scheduling is needed, async compacti
     on
07   * aborts. Sync compaction schedules.
08   * Either compaction type will also abort if a fatal signal is pending.
09   * In either case if the lock was locked, it is dropped and not regaine
     d.
10   *
11   * Returns true if compaction should abort due to fatal signal pending,
     or
12   *          async compaction due to need_resched()
13   * Returns false when compaction can continue (sync compaction might hav
     e
14   *          scheduled)
15   */
```

```
01  static bool compact_unlock_should_abort(spinlock_t *lock,
02                  unsigned long flags, bool *locked, struct compact_contro
     l *cc)
03  {
04          if (*locked) {
05                  spin_unlock_irqrestore(lock, flags);
06                  *locked = false;
07          }
08
```

```
09              if (fatal_signal_pending(current)) {
10                      cc->contended = true;
11                      return true;
12              }
13
14              if (need_resched()) {
15                      if (cc->mode == MIGRATE_ASYNC) {
16                              cc->contended = true;
17                              return true;
18                      }
19                      cond_resched();
20              }
21
22              return false;
23      }
```

If locked, unlock, and if there is an aborting element, set true to cc->contended and return true.

- Interrupting Factors
    - The SIGKILL signal is delayed.
    - During an asynchronous migration, there is a request for a lease schedule from a high-priority task.

## compact_trylock_irqsave()

mm/compaction.c

```
1  /*
2   * Compaction requires the taking of some coarse locks that are potentia
   lly
3   * very heavily contended. For async compaction, back out if the lock ca
   nnot
4   * be taken immediately. For sync compaction, spin on the lock if neede
   d.
5   *
6   * Returns true if the lock is held
7   * Returns false if the lock is not held and compaction should abort
8   */
```

```
01  static bool compact_trylock_irqsave(spinlock_t *lock, unsigned long *fla
   gs,
02                                              struct compact_control *
   cc)
03  {
04          if (cc->mode == MIGRATE_ASYNC) {
05                  if (!spin_trylock_irqsave(lock, *flags)) {
06                          cc->contended = true;
07                          return false;
08                  }
09          } else {
10                  spin_lock_irqsave(lock, *flags);
11          }
12
13          return true;
14  }
```

During asynchronous migration processing, if the spin lock used for the compaction competes with another CPU and the acquisition attempt fails at once, CC->contended is assigned true (compaction lock congestion) and returns false.

- In the case of synchronous migration, the lock is obtained unconditionally.

## mem_cgroup_page_lruvec()

mm/memcontrol.c

```
1   /**
2    * mem_cgroup_page_lruvec - return lruvec for isolating/putting an LRU p
     age
3    * @page: the page
4    * @zone: zone of the page
5    *
6    * This function is only safe when following the LRU page isolation
7    * and putback protocol: the LRU lock must be held, and the page must
8    * either be PageLRU() or the caller must have isolated/allocated it.
9    */

01  struct lruvec *mem_cgroup_page_lruvec(struct page *page, struct pglist_d
    ata *pgdat)
02  {
03          struct mem_cgroup_per_node *mz;
04          struct mem_cgroup *memcg;
05          struct lruvec *lruvec;
06
07          if (mem_cgroup_disabled()) {
08                  lruvec = &pgdat->lruvec;
09                  goto out;
10          }
11
12          memcg = page->mem_cgroup;
13          /*
14           * Swapcache readahead pages are added to the LRU - and
15           * possibly migrated - before they are charged.
16           */
17          if (!memcg)
18                  memcg = root_mem_cgroup;
19
20          mz = mem_cgroup_page_nodeinfo(memcg, page);
21          lruvec = &mz->lruvec;
22  out:
23          /*
24           * Since a node can be onlined after the mem_cgroup was created,
25           * we have to be prepared to initialize lruvec->zone here;
26           * and if offlined then reonlined, we need to reinitialize it.
27           */
28          if (unlikely(lruvec->pgdat != pgdat))
29                  lruvec->pgdat = pgdat;
30          return lruvec;
31  }
```

If memcg is enabled, it returns the lruvec of the node of memcg recorded on that page, and if memcg is disabled, it returns the lruvec of that node.


# Isolate LRU Page

## __isolate_lru_page()

mm/vmscan.c

```
01  /*
02   * Attempt to remove the specified page from its LRU.  Only take this pa
     ge
03   * if it is of the appropriate PageActive status.  Pages which are being
04   * freed elsewhere are also ignored.
```

```
05    *
06    * page:          page to consider
07    * mode:          one of the LRU isolation modes defined above
08    *
09    * returns 0 on success, -ve errno on failure.
10    */

01   int __isolate_lru_page(struct page *page, isolate_mode_t mode)
02   {
03           int ret = -EINVAL;
04
05           /* Only take pages on the LRU. */
06           if (!PageLRU(page))
07                   return ret;
08
09           /* Compaction should not handle unevictable pages but CMA can do
   so */
10           if (PageUnevictable(page) && !(mode & ISOLATE_UNEVICTABLE))
11                   return ret;
12
13           ret = -EBUSY;
14
15           /*
16            * To minimise LRU disruption, the caller can indicate that it o
   nly
17            * wants to isolate pages it will be able to operate on without
18            * blocking - clean pages for the most part.
19            *
20            * ISOLATE_ASYNC_MIGRATE is used to indicate that it only wants
   to pages
21            * that it is possible to migrate without blocking
22            */
23           if (mode & ISOLATE_ASYNC_MIGRATE) {
24                   /* All the caller can do on PageWriteback is block */
25                   if (PageWriteback(page))
26                           return ret;
27
28                   if (PageDirty(page)) {
29                           struct address_space *mapping;
30                           bool migrate_dirty;
31
32                           /*
33                            * Only pages without mappings or that have a
34                            * ->migratepage callback are possible to migrat
   e
35                            * without blocking. However, we can be racing w
   ith
36                            * truncation so it's necessary to lock the page
37                            * to stabilise the mapping as truncation holds
38                            * the page lock until after the page is removed
39                            * from the page cache.
40                            */
41                           if (!trylock_page(page))
42                                   return ret;
43
44                           mapping = page_mapping(page);
45                           migrate_dirty = !mapping || mapping->a_ops->migr
   atepage;
46                           unlock_page(page);
47                           if (!migrate_dirty)
48                                   return ret;
49                   }
50           }
51
52           if ((mode & ISOLATE_UNMAPPED) && page_mapped(page))
53                   return ret;
54
55           if (likely(get_page_unless_zero(page))) {
```

```
56                    /*
57                     * Be careful not to clear PageLRU until after we're
58                     * sure the page is not being freed elsewhere -- the
59                     * page release code relies on it.
60                     */
61                    ClearPageLRU(page);
62                    ret = 0;
63            }
64
65            return ret;
66    }
```

Attempt to remove LRU from the page. Returns 0 on success, -EINVAL or -EBUSY if it fails.

- In line 6~7 of the code, if it is not an LRU page, it will stop processing.
- In line 10~11 of the code, if the page is unevictable and the isolation mode does not support ISOLATE_UNEVICTABLE, the processing will be stopped.
- In line 23~26 of the code, if the page is operating in ISOLATE_ASYNC_MIGRATE mode, the page that is being WriteBack will be aborted with -EBUSY.
- If the migratepage handler function is not registered for a page set to Dirty in code lines 28~49, it will stop processing.
- When operating in ISOLATE_UNMAPPED mode on lines 52~53 of code, the mapping pages stop processing.
- In line 55~65 of the code, if the page is in use with a high probability, it clears the LRU flag and exits the function normally.

# Isolate movable page

The objects that use the following functions are used for the following purposes:

- Compaction
  - Migration of non-lru movable pages in direct-compliance
- Off-line Memory
  - Migration of non-LRU movable pages to off-line memory

## isolate_movable_page()

mm/migrate.c

```
01    int isolate_movable_page(struct page *page, isolate_mode_t mode)
02    {
03            struct address_space *mapping;
04
05            /*
06             * Avoid burning cycles with pages that are yet under __free_pag
       es(),
07             * or just got freed under us.
08             *
09             * In case we 'win' a race for a movable page being freed under
       us and
10             * raise its refcount preventing __free_pages() from doing its j
       ob
11             * the put_page() at the end of this block will take care of
12             * release this page, thus avoiding a nasty leakage.
13             */
```

```
14              if (unlikely(!get_page_unless_zero(page)))
15                      goto out;
16
17              /*
18               * Check PageMovable before holding a PG_lock because page's own
   er
19               * assumes anybody doesn't touch PG_lock of newly allocated page
20               * so unconditionally grapping the lock ruins page's owner side.
21               */
22              if (unlikely(!__PageMovable(page)))
23                      goto out_putpage;
24              /*
25               * As movable pages are not isolated from LRU lists, concurrent
26               * compaction threads can race against page migration functions
27               * as well as race against the releasing a page.
28               *
29               * In order to avoid having an already isolated movable page
30               * being (wrongly) re-isolated while it is under migration,
31               * or to avoid attempting to isolate pages being released,
32               * lets be sure we have the page lock
33               * before proceeding with the movable page isolation steps.
34               */
35              if (unlikely(!trylock_page(page)))
36                      goto out_putpage;
37
38              if (!PageMovable(page) || PageIsolated(page))
39                      goto out_no_isolated;
40
41              mapping = page_mapping(page);
42              VM_BUG_ON_PAGE(!mapping, page);
43
44              if (!mapping->a_ops->isolate_page(page, mode))
45                      goto out_no_isolated;
46
47              /* Driver shouldn't use PG_isolated bit of page->flags */
48              WARN_ON_ONCE(PageIsolated(page));
49              __SetPageIsolated(page);
50              unlock_page(page);
51
52              return 0;
53
54  out_no_isolated:
55              unlock_page(page);
56  out_putpage:
57              put_page(page);
58  out:
59              return -EBUSY;
60  }
```

isolates non-lru movable pages.

- In line 14~15 of the code, if the page has just become free, the function exits with the -EBUSY error.
- In line 22~23 of the code, if the page is not non-lru movable, the function exits with a -EBUSY error.
- In line 35~36 of the code, we get the page lock and check once again if the page is non-lru movable and if it is already isolate. If the page is not non-lru movable or has already been isolate, the function exits with a -EBUSY error.
- On lines 38~42 of code, if the page is not non-LRU movable or isolate, go to the out_no_isolated label.
    - If you call a function registered in the (*isolate_page) hook of a driver with a non-LRU movable applied (e.g. z3fold, zsmalloc, etc.) and it returns false, the function exits with an -

EBUSY error.

- Set the isolate flag in lines 46~49 of the code, and exit the function with a success (0) result.

# Isolatin Eligibility

## Isolation Eligibility

### Isolation for blocks?

Direct-Compaction, with the exception of sync-full, skips the isolation of that page block when it encounters a skip bit on the page block being scanned for fast compaction.

- If isolation has recently failed in the page block, skip the block.

The following illustration shows the use of the usemap's skip bit for page blocks.

 (http://jake.dothome.co.kr/wp-content/uploads/2016/07/compaction-5.png)

### Isolation Force

The following performance conditions force an attempt to isolate all blocks by setting a ignore_skip_hint internally, regardless of whether the skip bit is set or not.

- When using direct-compact using compact mode except for sync-full, which has the highest priority
- When using manual-compact
- With kcompactd

### Mobility characteristics of Pageblock

As shown in the following figure, the type of the page block has a representative mobility characteristic with more than half (50%) of the page (migratetype) type in the block.

- When using direct-compaction, which should be fast, the compaction is performed only on movable page blocks.



(http://jake.dothome.co.kr/wp-content/uploads/2016/07/pageblock-1.png)

## isolation_suitable()

mm/compaction.c

```
1  /* Returns true if the pageblock should be scanned for pages to isolate. */
2  static inline bool isolation_suitable(struct compact_control *cc,
3                                        struct page *page)
4  {
5          if (cc->ignore_skip_hint)
6                  return true;
7
8          return !get_pageblock_skip(page);
9  }
```

Check whether it is okay to try isolation from the block.

- If the ignore_skip_hint is set in line 5~6 of the code, it returns true so that the page block can proceed with isolation regardless of whether it is skipped or not.
- Line 8 of the code returns false to skip if isolation has failed in that page block recently.
    - usemap stores the skip bit of each page block.

The following figure shows the processing of the isolation_suitable() function, which returns the isolation of the block.



(http://jake.dothome.co.kr/wp-content/uploads/2016/07/isolation_suitable-1b.png)

# Suitability as a migration source

## suitable_migration_source()

mm/compaction.c

```
01  static bool suitable_migration_source(struct compact_control *cc,
02                                              struct page *page)
03  {
04          int block_mt;
05
06          if ((cc->mode != MIGRATE_ASYNC) || !cc->direct_compaction)
07                  return true;
08
09          block_mt = get_pageblock_migratetype(page);
10
11          if (cc->migratetype == MIGRATE_MOVABLE)
12                  return is_migrate_movable(block_mt);
13          else
14                  return block_mt == cc->migratetype;
15  }
```

Returns whether a page contained in a page block is eligible for the migration target. In order to make direct-compaction in asynchronous mode lighter and faster, it returns true only if the block migrate type is the same as the requested migrate type. Otherwise, it always returns true.

- In line 6~7 of code, it returns true to unconditionally isolate the block and handle the migrated pages in the following 2 cases.
    - Non-migrate_async requests
    - Requesting manual compaction or kcompactd that are not of type direct_compaction
- When requesting a movable assignment in line 9~12 of the code, it returns whether the movable and CMA block types are present.
    - Unmovable, reclimable, highatomic, and iolate types cannot be subject to migration.
- In line 13~14 of the code, if the movable allocation request is not the same, it returns whether the block type corresponding to the requested type is the same.

The following figure shows the processing of the suitable_migration_source() function, which returns whether the block is appropriate as a migration source.



(http://jake.dothome.co.kr/wp-content/uploads/2016/07/suitable_migration_source-1a.png)

# Suitability as a migration target
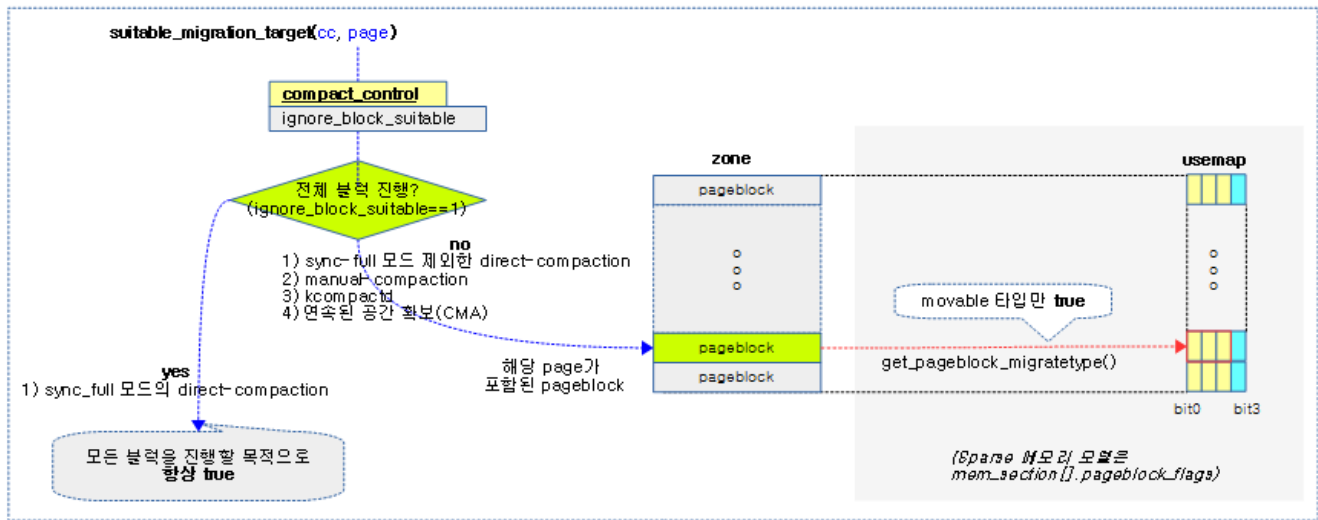
## suitable_migration_target()

mm/compaction.c

```
01  /* Returns true if the page is within a block suitable for migration to
    */
02  static bool suitable_migration_target(struct compact_control *cc,
03                                                  struct page *pag
    e)
04  {
05          /* If the page is a large free page, then disallow migration */
06          if (PageBuddy(page)) {
07                  /*
08                   * We are checking page_order without zone->lock taken.
    But
09                   * the only small danger is that we skip a potentially s
    uitable
10                   * pageblock, so it's not worth to check order for valid
    range.
11                   */
12                  if (page_order_unsafe(page) >= pageblock_order)
13                          return false;
14          }
15
16          if (cc->ignore_block_suitable)
17                  return true;
18
19          /* If the block is MIGRATE_MOVABLE or MIGRATE_CMA, allow migrati
    on */
20          if (is_migrate_movable(get_pageblock_migratetype(page)))
21                  return true;
22
23          /* Otherwise skip the block */
24          return false;
25  }
```

Returns whether a page contained in a page block is eligible for the migration target. However, direct-compaction in sync-full mode always returns true.

- In code lines 6~14, a free page on a buddy system that is larger than a page block returns false.
- In lines 16~17 of the code, return true to unconditionally target regardless of the type of page block type.
  - The ignore_block_suitable is only set when direct-component in sync-full mode is enabled.
- In lines 20~24 of code, it returns true only if the page block is of type movable, and false otherwise.

The following figure shows the processing of the suitable_migration_target() function, which returns whether the block is appropriate as a migration target.

(http://jake.dothome.co.kr/wp-content/uploads/2016/07/suitable_migration_target-1.png)

# Specifying a Page Block Skip

## update_pageblock_skip()

mm/compaction.c

```
1   /*
2    * If no pages were isolated then mark this pageblock to be skipped in the
3    * future. The information is later cleared by __reset_isolation_suitable().
4    */
```
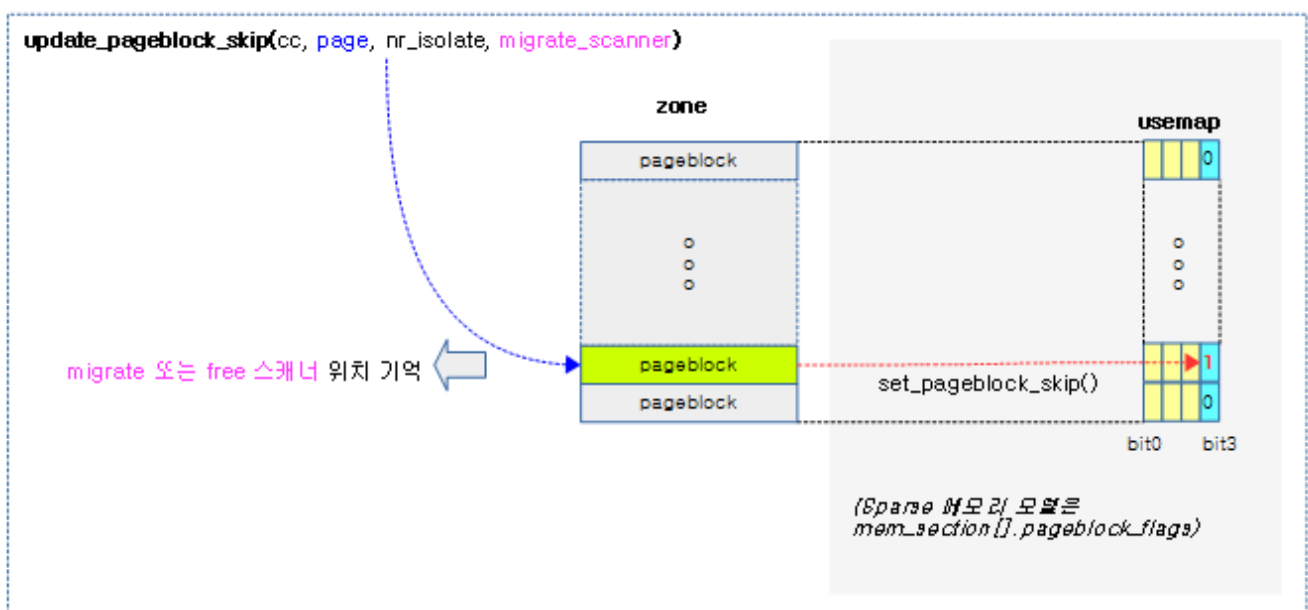
```
01  static void update_pageblock_skip(struct compact_control *cc,
02                      struct page *page, unsigned long nr_isolated,
03                      bool migrate_scanner)
04  {
05          struct zone *zone = cc->zone;
06          unsigned long pfn;
07
08          if (cc->no_set_skip_hint)
09                  return;
10
11          if (!page)
12                  return;
13
14          if (nr_isolated)
15                  return;
16
17          set_pageblock_skip(page);
18
19          pfn = page_to_pfn(page);
20
21          /* Update where async and sync compaction should restart */
22          if (migrate_scanner) {
23                  if (pfn > zone->compact_cached_migrate_pfn[0])
24                          zone->compact_cached_migrate_pfn[0] = pfn;
25                  if (cc->mode != MIGRATE_ASYNC &&
26                      pfn > zone->compact_cached_migrate_pfn[1])
27                          zone->compact_cached_migrate_pfn[1] = pfn;
28          } else {
29                  if (pfn < zone->compact_cached_free_pfn)
30                          zone->compact_cached_free_pfn = pfn;
31          }
32  }
```

If there are no isolated pages processed, set the migrate skip bit in the page block and set the start pfn of the migrate or free scanner to that page.

- In line 8~15 of the code, the following 3 cases abort to avoid setting the skip bit of the corresponding page block.
  - When no_set_skip_hint is set
    - It is set when isolation is in progress to secure the CMA area.
  - If you don't have a page
  - If you don't have an isolated page
- In line 17 of code, set the skip bit for that page block.
- If a request comes from the migrate scanner routine on lines 22~27 of code, update the pfn to start the migrate scan.
- If a request comes from the free scanner routine on code lines 28~31, update the free scan start pfn.

The following figure shows the handling of the update_pageblock_skip() function, which sets the skip bit to prohibit isolation of the block.



(http://jake.dothome.co.kr/wp-content/uploads/2016/07/update_pageblock_skip-1.png)

# Reset Zone Skip Bit

## __reset_isolation_suitable()

mm/compaction.c

```
1  /*
2   * This function is called to clear all cached information on pageblocks
     that
3   * should be skipped for page isolation when the migrate and free page s
     canner
4   * meet.
5   */
01 static void __reset_isolation_suitable(struct zone *zone)
```

```
02  {
03          unsigned long start_pfn = zone->zone_start_pfn;
04          unsigned long end_pfn = zone_end_pfn(zone);
05          unsigned long pfn;
06
07          zone->compact_blockskip_flush = false;
08
09          /* Walk the zone and mark every pageblock as suitable for isolat
    ion */
10          for (pfn = start_pfn; pfn < end_pfn; pfn += pageblock_nr_pages)
    {
11                  struct page *page;
12
13                  cond_resched();
14
15                  page = pfn_to_online_page(pfn);
16                  if (!page)
17                          continue;
18                  if (zone != page_zone(page))
19                          continue;
20                  if (pageblock_skip_persistent(page))
21                          continue;
22
23                  clear_pageblock_skip(page);
24          }
25
26          reset_cached_positions(zone);
27  }
```

Reset the isolate and free scan start address of the zone, and clear all skip bits in the page block. Excludes the skip bit for compound pages that use more than a page block order.

- The migration target area is the start pfn and the end pfn of the zone.
  - Assign the starting pfn of the zone to compact_cached_migrate_pfn[0..1] as the pfn that the isolate scanner will start.
  - Assign the end pfn of the zone to the compact_cached_free_pfn with the pfn that the free scanner will start.
- pageblock->flags
  - If not the Sparse memory model, zone->pageblock_flags specifies usemap
  - If it's a Sparse memory model, the usemap specified by mem_section[].pageblock_flags

The following figure shows the process of the __reset_isolation_suitable() function, which clears the skip bit to allow isolation of all blocks in the zone, and resets the starting position of the scanners.

(http://jake.dothome.co.kr/wp-content/uploads/2016/07/reset_isolation_suitable-1.png)

# Isolation Freepages

The following functions are used for compaction when the high order page is insufficient, where direct-compact, manual-compact, and kcompactd functions are used.

## isolate_freepages()

mm/compaction.c

```
1  /*
2   * Based on information in the current compact_control, find blocks
3   * suitable for isolating free pages from and then isolate them.
4   */
01 static void isolate_freepages(struct compact_control *cc)
02 {
03         struct zone *zone = cc->zone;
04         struct page *page;
05         unsigned long block_start_pfn;  /* start of current pageblock */
06         unsigned long isolate_start_pfn; /* exact pfn we start at */
07         unsigned long block_end_pfn;    /* end of current pageblock */
08         unsigned long low_pfn;        /* lowest pfn scanner is able to scan */
09         struct list_head *freelist = &cc->freepages;
10
11         /*
12          * Initialise the free scanner. The starting point is where we last
13          * successfully isolated from, zone-cached value, or the end of the
14          * zone when isolating for the first time. For looping we also need
15          * this pfn aligned down to the pageblock boundary, because we do
16          * block_start_pfn -= pageblock_nr_pages in the for loop.
17          * For ending point, take care when isolating in last pageblock of a
18          * a zone which ends in the middle of a pageblock.
19          * The low boundary is the end of the pageblock the migration scanner
```

```
 20                 * is using.
 21                 */
 22             isolate_start_pfn = cc->free_pfn;
 23             block_start_pfn = pageblock_start_pfn(cc->free_pfn);
 24             block_end_pfn = min(block_start_pfn + pageblock_nr_pages,
 25                                             zone_end_pfn(zone));
 26             low_pfn = pageblock_end_pfn(cc->migrate_pfn);
 27
 28             /*
 29              * Isolate free pages until enough are available to migrate the
 30              * pages on cc->migratepages. We stop searching if the migrate
 31              * and free page scanners meet or enough free pages are isolate
    d.
 32              */
 33             for (; block_start_pfn >= low_pfn;
 34                             block_end_pfn = block_start_pfn,
 35                             block_start_pfn -= pageblock_nr_pages,
 36                             isolate_start_pfn = block_start_pfn) {
 37                     /*
 38                      * This can iterate a massively long zone without findin
    g any
 39                      * suitable migration targets, so periodically check if
    we need
 40                      * to schedule, or even abort async compaction.
 41                      */
 42                     if (!(block_start_pfn % (SWAP_CLUSTER_MAX * pageblock_nr
    _pages))
 43                                             && compact_should_abort
    (cc))
 44                             break;
 45
 46                     page = pageblock_pfn_to_page(block_start_pfn, block_end_
    pfn,
 47                                                                         z
    one);
 48                     if (!page)
 49                             continue;
 50
 51                     /* Check the block is suitable for migration */
 52                     if (!suitable_migration_target(cc, page))
 53                             continue;
 54
 55                     /* If isolation recently failed, do not retry */
 56                     if (!isolation_suitable(cc, page))
 57                             continue;
 58
 59                     /* Found a block suitable for isolating free pages from.
     */
 60                     isolate_freepages_block(cc, &isolate_start_pfn, block_en
    d_pfn,
 61                                             freelist, false);
 62
 63                     /*
 64                      * If we isolated enough freepages, or aborted due to lo
    ck
 65                      * contention, terminate.
 66                      */
 67                     if ((cc->nr_freepages >= cc->nr_migratepages)
 68                                             || cc->contende
    d) {
 69                             if (isolate_start_pfn >= block_end_pfn) {
 70                                     /*
 71                                      * Restart at previous pageblock if more
 72                                      * freepages can be isolated next time.
 73                                      */
 74                                     isolate_start_pfn =
 75                                             block_start_pfn - pageblock_nr_p
    ages;
```

```
76                              }
77                              break;
78                      } else if (isolate_start_pfn < block_end_pfn) {
79                              /*
80                               * If isolation failed early, do not continue
81                               * needlessly.
82                               */
83                              break;
84                      }
85              }
86
87              /* __isolate_free_page() does not map the pages */
88              map_pages(freelist);
89
90              /*
91               * Record where the free scanner will restart next time. Either
     we
92               * broke from the loop and set isolate_start_pfn based on the la
     st
93               * call to isolate_freepages_block(), or we met the migration sc
     anner
94               * and the loop terminated due to isolate_start_pfn < low_pfn
95               */
96              cc->free_pfn = isolate_start_pfn;
97      }
```

- Start with cc->free_pfn which was discontinued for the start pfn of the free scanner in code lines 22~25. Then specify the start and end pfn of the block.
  - The free scanner scans the block by moving it down the zone. When isolating free pages within a page block, proceed in the upward direction.
- In order to prevent the block that the migrate scanner is working on at line 26, substitute the end pfn of the page block where the migrate scanner is active in the low_pfn.
- In lines 33~36 of code, the free scanner tverses in the downward direction where the migrate scanner is located, decreasing by page blocks.
- In line 42~44 of code, if you isolate the zones all at once with a free scanner, the area is too large and it is checked for compaction interrupts once in the middle.
  - Check Frequency
    - SWAP_CLUSTER_MAX (32) pages are checked in blocks.
  - Interrupting Factors
    - Asynchronous compaction processing is in progress and there is a preemption request from a higher priority task.
- In lines 46~49 of the code, get the first page of the page block. If the pfn scope is not in the requested zone, skip it.
  - If the last page block of the zone is partial, it should also be skipped
- In lines 52~53 of code, if the page block is not suitable as a target for migration, skip it.
- In line 56~57 of code, if the page block is not suitable for isolation, skip it.
  - Returns false to skip if isolation has recently been canceled from the pageblock.
- In lines 60~61 of the code, isolate the page block and move it to the cc->freepages list.
- In line 67~77 of the code, if there are more free pages secured by the free scanner than migratable pages obtained by the migrate scanner, or if there is a lock contention situation, the loop will be aborted.
- If a single block is not scanned in code lines 78~84, the loop is aborted.

- For a list of free pages isolated from code line 88, remove each page from the LRU list, and decompose it into a 0 order.
- In line 96 of code, note the location of the free scanner for the next scan.

---

# Isolate freepages block

The following function is called to be used for two purposes:

- Compaction
  - isolate_freepages() function
    - It is used for direct-compact, manual-compact, and kcompactd for compaction to obtain high order pages.
- CMA
  - isolate_freepages_range() function.
    - It is used to clear the requested range of the CMA area.

## isolate_freepages_block()

mm/compaction.c -1/2-

```
 1  /*
 2   * Isolate free pages onto a private freelist. If @strict is true, will
      abort
 3   * returning 0 on any invalid PFNs or non-free pages inside of the pageb
      lock
 4   * (even though it may still end up isolating some pages).
 5   */

01  static unsigned long isolate_freepages_block(struct compact_control *cc,
02                                      unsigned long *start_pfn,
03                                      unsigned long end_pfn,
04                                      struct list_head *freelist,
05                                      bool strict)
06  {
07          int nr_scanned = 0, total_isolated = 0;
08          struct page *cursor, *valid_page = NULL;
09          unsigned long flags = 0;
10          bool locked = false;
11          unsigned long blockpfn = *start_pfn;
12          unsigned int order;
13
14          cursor = pfn_to_page(blockpfn);
15
16          /* Isolate free pages. */
17          for (; blockpfn < end_pfn; blockpfn++, cursor++) {
18                  int isolated;
19                  struct page *page = cursor;
20
21                  /*
22                   * Periodically drop the lock (if held) regardless of it
      s
23                   * contention, to give chance to IRQs. Abort if fatal si
      gnal
24                   * pending or async compaction detects need_resched()
25                   */
26                  if (!(blockpfn % SWAP_CLUSTER_MAX)
```

```
27                         && compact_unlock_should_abort(&cc->zone->lock, flag
s,
28                                                         &locked,
cc))
29                     break;
30
31              nr_scanned++;
32              if (!pfn_valid_within(blockpfn))
33                     goto isolate_fail;
34
35              if (!valid_page)
36                     valid_page = page;
37
38              /*
39               * For compound pages such as THP and hugetlbfs, we can
save
40               * potentially a lot of iterations if we skip them at on
ce.
41               * The check is racy, but we can consider only valid val
ues
42               * and the only danger is skipping too much.
43               */
44              if (PageCompound(page)) {
45                     const unsigned int order = compound_order(page);
46
47                     if (likely(order < MAX_ORDER)) {
48                             blockpfn += (1UL << order) - 1;
49                             cursor += (1UL << order) - 1;
50                     }
51                     goto isolate_fail;
52              }
53
54              if (!PageBuddy(page))
55                     goto isolate_fail;
56
57              /*
58               * If we already hold the lock, we can skip some recheck
ing.
59               * Note that if we hold the lock now, checked_pageblock
was
60               * already set in some previous iteration (or strict is
true),
61               * so it is correct to skip the suitable migration targe
t
62               * recheck as well.
63               */
64              if (!locked) {
65                     /*
66                      * The zone lock must be held to isolate freepag
es.
67                      * Unfortunately this is a very coarse lock and
can be
68                      * heavily contended if there are parallel alloc
ations
69                      * or parallel compactions. For async compaction
do not
70                      * spin on the lock and we acquire the lock as l
ate as
71                      * possible.
72                      */
73                     locked = compact_trylock_irqsave(&cc->zone->loc
k,
74                                                         &flags,
cc);
75                     if (!locked)
76                             break;
77
78                     /* Recheck this is a buddy page under lock */
```

```
79                         if (!PageBuddy(page))
80                                 goto isolate_fail;
81                 }
```

- In code lines 14~19, traverse the page from the beginning pfn to the end pfn within the block.
- In line 26~29 of code, I have a large scope to process while scanning a block, so I periodically unlock it if it is locked, and if there is an interrupting element, I set ture to cc->contended and exit the loop.
    - Check Frequency
        - SWAP_CLUSTER_MAX (32) pages per block
    - Interrupting Elements
        - The SIGKILL signal is delayed.
        - If there is a preemption request from a higher priority task during an asynchronous migration.
- Increment the scan counter at line 31 of code.
- On lines 32~36 of the code, if it is not a valid page, go to the isolate_fail label. If it is the first valid page, remember it.
- On lines 44~52 of the code, go to the isolate_fail label if it is a compound page.
- In code lines 54~55, go to the isolate_fail label if the page is not free managed by the buddy system.
- In code lines 64~81, if there is no lock, get a lock, check whether the buddy system is managing the free page again, and if not, move to the isolate_fail label.

mm/compaction.c -2/2-

```
01                         /* Found a free page, will break it into order-0 pages
   */
02                         order = page_order(page);
03                         isolated = __isolate_free_page(page, order);
04                         if (!isolated)
05                                 break;
06                         set_page_private(page, order);
07
08                         total_isolated += isolated;
09                         cc->nr_freepages += isolated;
10                         list_add_tail(&page->lru, freelist);
11
12                         if (!strict && cc->nr_migratepages <= cc->nr_freepages)
   {
13                                 blockpfn += isolated;
14                                 break;
15                         }
16                         /* Advance to the end of split page */
17                         blockpfn += isolated - 1;
18                         cursor += isolated - 1;
19                         continue;
20
21  isolate_fail:
22                         if (strict)
23                                 break;
24                         else
25                                 continue;
26
27                 }
28
29         if (locked)
```

```
30                        spin_unlock_irqrestore(&cc->zone->lock, flags);
31
32                /*
33                 * There is a tiny chance that we have read bogus compound_order
      (),
34                 * so be careful to not go outside of the pageblock.
35                 */
36                if (unlikely(blockpfn > end_pfn))
37                        blockpfn = end_pfn;
38
39                trace_mm_compaction_isolate_freepages(*start_pfn, blockpfn,
40                                        nr_scanned, total_isolated);
41
42                /* Record how far we have got within the block */
43                *start_pfn = blockpfn;
44
45                /*
46                 * If strict isolation is requested by CMA then check that all t
      he
47                 * pages requested were isolated. If there were any failures, 0
      is
48                 * returned and CMA will fail.
49                 */
50                if (strict && blockpfn < end_pfn)
51                        total_isolated = 0;
52
53                /* Update the pageblock-skip if the whole pageblock was scanned
      */
54                if (blockpfn == end_pfn)
55                        update_pageblock_skip(cc, valid_page, total_isolated, fa
      lse);
56
57                cc->total_free_scanned += nr_scanned;
58                if (total_isolated)
59                        count_compact_events(COMPACTISOLATED, total_isolated);
60                return total_isolated;
61        }
```

- In code lines 2~6, separate the free page from the free list of the buddy system and record the order.
- In line 8~9 of the code, increment the isolate counter and the secured free page counter.
- Add the page you separated from line 10 to the list of free scanners.
- If the @strict is true in line 12~15 of the code, it is requested to obtain contiguous pages in the CMA area, and if any of them fail, the processing will be aborted. If the @strict is zero and the free scanner has more pages than the migration scanner has obtained, the process will still be stopped.
- On lines 17~19 of the code, let the loop continue for the next page.
- In code lines 21~25, isolate_fail: Label. If the @strict is 1, it will exit the loop immediately if it fails, otherwise it will continue the loop.
- Unlock the lock obtained in code lines 29~30.
- Limit the ongoing pfn in line 36~37 to not exceed the ending pfn.
- Update the @start_pfn value received as an input/output argument in line 43 with the current pfn value.
- In code lines 50~51, if the @strict is true, it is requested to obtain contiguous pages in the CMA area, in which case the total_isolated value is changed to 0 to return the failure value to 0.
- I went all the way to the end of line 54~55 of the code, but it failed, so I marked this page block with a skip mark.

- At line 57 of code, add the number of pages scanned and the number of pages isolated, and return the number of isolated pages.

## Isolate Free Page

### __isolate_free_page()

mm/page_alloc.c

```
01   int __isolate_free_page(struct page *page, unsigned int order)
02   {
03           unsigned long watermark;
04           struct zone *zone;
05           int mt;
06
07           BUG_ON(!PageBuddy(page));
08
09           zone = page_zone(page);
10           mt = get_pageblock_migratetype(page);
11
12           if (!is_migrate_isolate(mt)) {
13                   /*
14                    * Obey watermarks as if the page was being allocated. We can
15                    * emulate a high-order watermark check with a raised order-0
16                    * watermark, because we already know our high-order page
17                    * exists.
18                    */
19                   watermark = min_wmark_pages(zone) + (1UL << order);
20                   if (!zone_watermark_ok(zone, 0, watermark, 0, ALLOC_CMA))
21                           return 0;
22
23                   __mod_zone_freepage_state(zone, -(1UL << order), mt);
24           }
25
26           /* Remove page from free list */
27           list_del(&page->lru);
28           zone->free_area[order].nr_free--;
29           rmv_page_order(page);
30
31           /*
32            * Set the pageblock if the isolated page is at least half of a
33            * pageblock
34            */
35           if (order >= pageblock_order - 1) {
36                   struct page *endpage = page + (1 << order) - 1;
37                   for (; page < endpage; page += pageblock_nr_pages) {
38                           int mt = get_pageblock_migratetype(page);
39                           if (!is_migrate_isolate(mt) && !is_migrate_cma(mt)
40                               && !is_migrate_highatomic(mt))
41                                   set_pageblock_migratetype(page,
42                                                   MIGRATE_MOVABLE);
43                   }
44           }
45
46
47           return 1UL << order;
48   }
```

If the order page is a free page, it will be decomposed and returned as order-0 free page. If this is normal, the number of pages corresponding to 2^order is returned.

- In line 9~10 of the code, we know the zone and migrate types that correspond to the page.
- In line 12~24 of the code, if the migrate type is impossible to isolate, it will request order 0 and return 0 as a failure if it does not pass the watermark boundary (OK) even with the value of low watermark + order number of pages. If it can pass, reduce the number of free pages by the number of order pages.
- In line 27~29 of the code, remove the buddy flag and clear the private value of 0 where the order information is recorded in order to separate the page from the buddy system and remove the buddy information from the page.
- When freeing the order page in line 35~44, if the page block is more than 50%, the type of the page block can be changed to movable if the existing type of the page block is unmovable and reclaimable. However, the isolate, CMA, and highatomic types cannot be changed.
    - Since there can be multiple page blocks within an order page, it traverses page block by page to the end of the order page, and if the page block is not all of the three types: isolate, cma, and highatomic, it is changed to movable.

## split_page()

mm/page_alloc.c

```
 1  /*
 2   * split_page takes a non-compound higher-order page, and splits it into
 3   * n (1<<order) sub-pages: page[0..n]
 4   * Each sub-page must be freed individually.
 5   *
 6   * Note: this is probably too low level an operation for use in drivers.
 7   * Please consult with lkml before using this in your driver.
 8   */

01  void split_page(struct page *page, unsigned int order)
02  {
03          int i;
04
05          VM_BUG_ON_PAGE(PageCompound(page), page);
06          VM_BUG_ON_PAGE(!page_count(page), page);
07
08          for (i = 1; i < (1 << order); i++)
09                  set_page_refcounted(page + i);
10          split_page_owner(page, order);
11  }
12  EXPORT_SYMBOL_GPL(split_page);
```

Set the reference counter to 1 on all split pages. Then, move all the owners of the existing pages.

# Get new page hook function

## compaction_alloc()

mm/compaction.c

```
 1  /*
```

```
 2    * This is a migrate-callback that "allocates" freepages by taking pages
 3    * from the isolated freelists in the block we are migrating to.
 4    */

01   static struct page *compaction_alloc(struct page *migratepage,
02                                         unsigned long data,
03                                         int **result)
04   {
05          struct compact_control *cc = (struct compact_control *)data;
06          struct page *freepage;
07
08          /*
09           * Isolate free pages if necessary, and if we are not aborting d
     ue to
10           * contention.
11           */
12          if (list_empty(&cc->freepages)) {
13                  if (!cc->contended)
14                          isolate_freepages(cc);
15
16                  if (list_empty(&cc->freepages))
17                          return NULL;
18          }
19
20          freepage = list_entry(cc->freepages.next, struct page, lru);
21          list_del(&freepage->lru);
22          cc->nr_freepages--;
23
24          return freepage;
25   }
```

Returns the leading free page from the cc->freepages list isolated from the free scanner. If there are no
free pages to return, run the free scanner. If it fails, it returns null.


## Put new page hook function

### compaction_free()

mm/compaction.c

```
1   /*
2    * This is a migrate-callback that "frees" freepages back to the isolate
    d
3    * freelist.  All pages on the freelist are from the same zone, so there
    is no
4    * special handling needed for NUMA.
5    */

1   static void compaction_free(struct page *page, unsigned long data)
2   {
3          struct compact_control *cc = (struct compact_control *)data;
4
5          list_add(&page->lru, &cc->freepages);
6          cc->nr_freepages++;
7   }
```

Add the page back to the cc->freepages list.

# consultation

---

## 4 thoughts to "Zoned Allocator -8- (Direct Compact-Isolation)"

**KWON YONGBEOM**

2020-12-26 23:22 (http://jake.dothome.co.kr/zonned-allocator-isolation/#comment-297399)

Report what appears to be a typo.
In the __isolate_free_page() code description, there are some parts
such as "from the corresponding list of the buddy system on code lines 27~29" and the
code line does not match the description.
I appreciate it

RESPONSE (/ZONNED-ALLOCATOR-ISOLATION/?REPLYTOCOM=297399#RESPOND)

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

2020-12-27 19:35 (http://jake.dothome.co.kr/zonned-allocator-isolation/#comment-297552)

Thanks for looking for the typo.

The blog code was in version 4.0 and the description was in version 5.4.
The code
has been changed to 5.4 code, and the explanations in lines 35~44 have been slightly enhanced.

I would be grateful if you could continue with the kernel analysis and find any oddities.

I appreciate it.

RESPONSE (/ZONNED-ALLOCATOR-ISOLATION/?REPLYTOCOM=297552#RESPOND)

**KWON YONGBEOM**

2020-12-28 17:01 (http://jake.dothome.co.kr/zonned-allocator-isolation/#comment-297753)

yes ^^ and I think the kernel versions are v4.x, v5.0, v5.4 and so on.
I wondered if the __reset_isolation_suitable() function had also been changed (as of v5.1) and not updated, but
it says kernel v5.0 at the top of the page.

RESPONSE (/ZONNED-ALLOCATOR-ISOLATION/?REPLYTOCOM=297753#RESPOND)

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

2020-12-28 17:22 (http://jake.dothome.co.kr/zonned-allocator-isolation/#comment-297757)

Yes, that's right. I wrote in 4.0 in the early days of my blog, and
I've been switching to 5.x since last year.
Most of the memory was used in the 5.0 source code, and then the scheduler was upgraded to the LTS version of 5.4.

Sometimes, however, I find a source that I haven't changed. ^^;
The __reset_isolation_suitable() function has also been modified.

I appreciate it.

응답 (/ZONNED-ALLOCATOR-ISOLATION/?REPLYTOCOM=297757#RESPOND)

# 댓글 남기기

이메일은 공개되지 않습니다. 필수 입력창은 * 로 표시되어 있습니다

댓글

이름 *

이메일 *

웹사이트

댓글 작성

❮ Zoned Allocator -7- (Direct Compact) (http://jake.dothome.co.kr/zonned-allocator-compaction/)

Zoned Allocator -9- (Direct Compact-Migration) ❯ (http://jake.dothome.co.kr/zonned-allocator-migration/)

문c 블로그 (2015 ~ 2023)