# Slub Memory Allocator -6- (Assign Object)

📅 2016-06-06 (http://jake.dothome.co.kr/slub-object-alloc/)　　👤 Moon Young-il (http://jake.dothome.co.kr/author/admin/)　　📁 Linux Kernel (http://jake.dothome.co.kr/category/linux/)

<kernel v5.0>

## Slub Object Assignment (Fastpath, Slowpath)

Slub object assignments consist of four levels of slowpath assignments, starting with the fastest and the slowest at most.

**Fastpath**

Fastpath action atomically fetches the first object from the freelist of the per cpu slab cache. If there is no slab object to assign, move on to the Slowpath action.

- FastPath Steps:
    - Assign a free object on the fly from c->freelist
    - For reference, the abbreviation means:
        - s: kmem_cache struct pointing to the slab cache
        - C: S->cpu_slab of Per-CPU (kmem_cache_cpu structure)
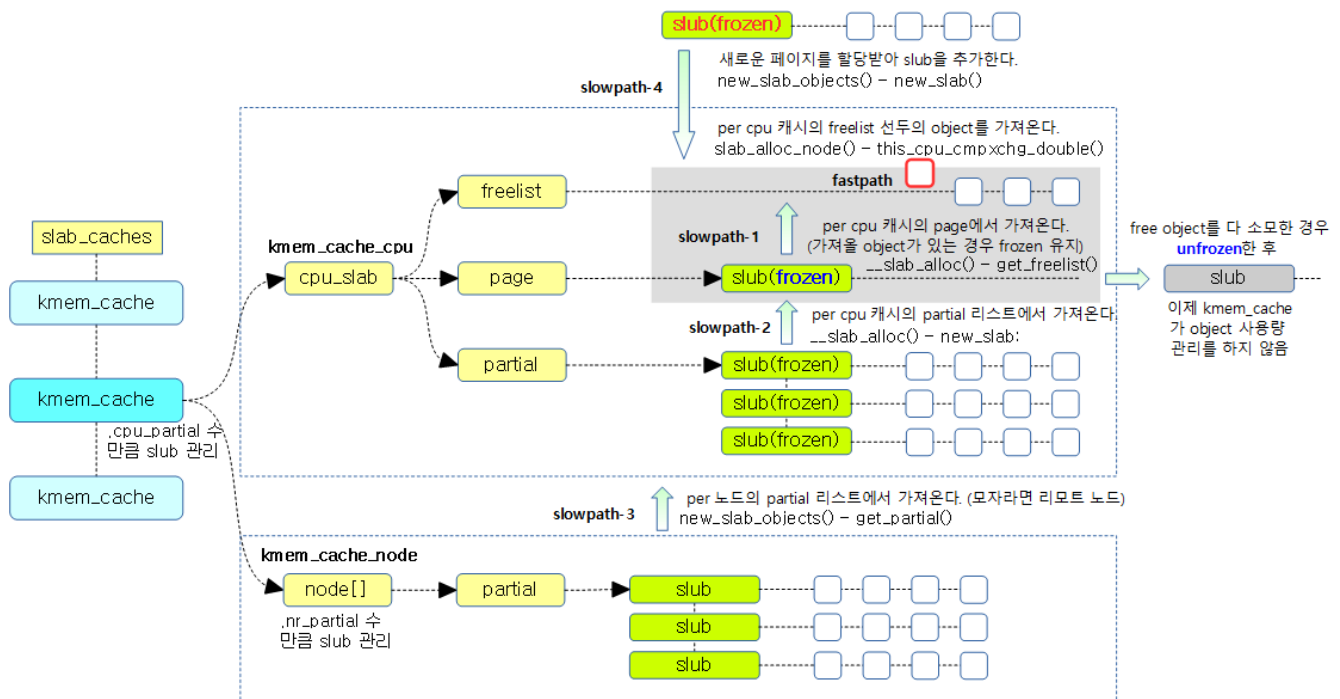        - n: s->node[] in per-node (kmem_cache_node struct)

**Slowpath**

If c->freelist is empty and there is no object to retrieve, proceed with the slowpath operation as shown below, refill it, and then retry the fastpath operation.

- Slowpath Step 1:
    - c->page->freelist —(이동)—> c->freelist
    - Objects from the page->freelist of the per CPU cache are migrated to the freelist of the per CPU cache. If there is a free object, freeze the page. If there is no free object, proceed to step 2.
- Slowpath Step 2:
    - c->partial —(이동)—> c->page
    - Migrate the first slap page from the partial list of the per cpu cache to the page of the per cpu cache. If the partial list of per CPU cache is empty, proceed to step 3.
- Slowpath Step 3:
    - n->partial—(moves first slap page)—> c->page
    - n->partial—(moved a few places from the second slab page)—> c->partial
    - First, migrate the first slap page from the partial list of per nodes to a page in the per CPU cache.

- Second, starting with the next slap page in the partial list of per nodes, add it to the partial list of the per CPU cache.
    - However, if the number of free objects on the slab page exceeds half of the s->cpu_partial, only the slab page is migrated and the repeat is no more.
  - If the partial list of per nodes is empty, it will retry from the partial list of other remote nodes. If even the other remote per nodes are unavailable (limiting the remote node's memory usage beyond a certain percentage), proceed to step 4.
- Slowpath Step 4:
  - Page allocator (buddy system)—(assign)—> c->page & c->freelist
  - After receiving the page assignment from the page allocator (buddy system) and initializing it as a slab page, it directly connects to c->page and connects the first free object to c->freelist.

The following illustration shows the priority of the slab object being assigned.



(http://jake.dothome.co.kr/wp-content/uploads/2016/06/slub-1e.png)

# Frozened Slub Page

The SLUB allocator manages the SLUB page from the per-cpu for fast allocation of SLUB objects without locks.

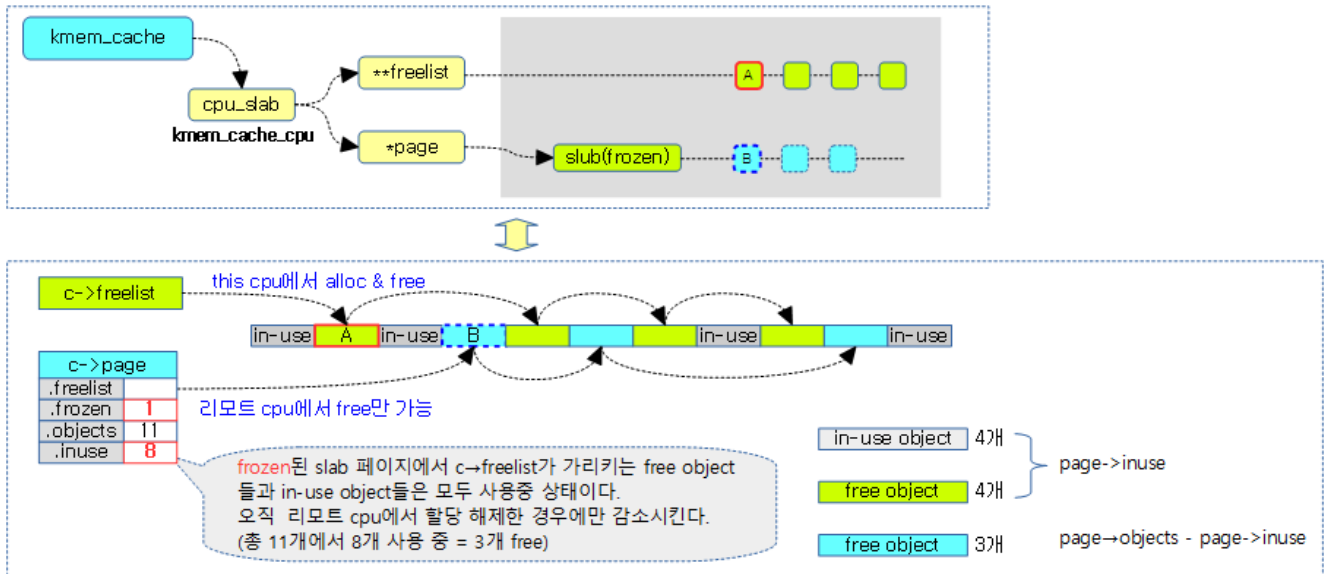- cmpxchg double atomic operation to further delete the list.

If page->frozen is set to 1, the freelist of the slab page is locked so that the current task manages it directly.

## Manage two freelists

In fact, the management of the freelist is divided into two, and it works as follows.

- c->freelist
  - Currently, the CPU manages the freelist, so it is possible to allocate and deallocate slub objects (free).
- page->freelist
  - Other CPUs can only be freed here.

The following figure shows how the slub page is operated with 11 objects and fozen, and the freelist is managed in two.



(http://jake.dothome.co.kr/wp-content/uploads/2016/06/frozen-slub-page-1c.png)
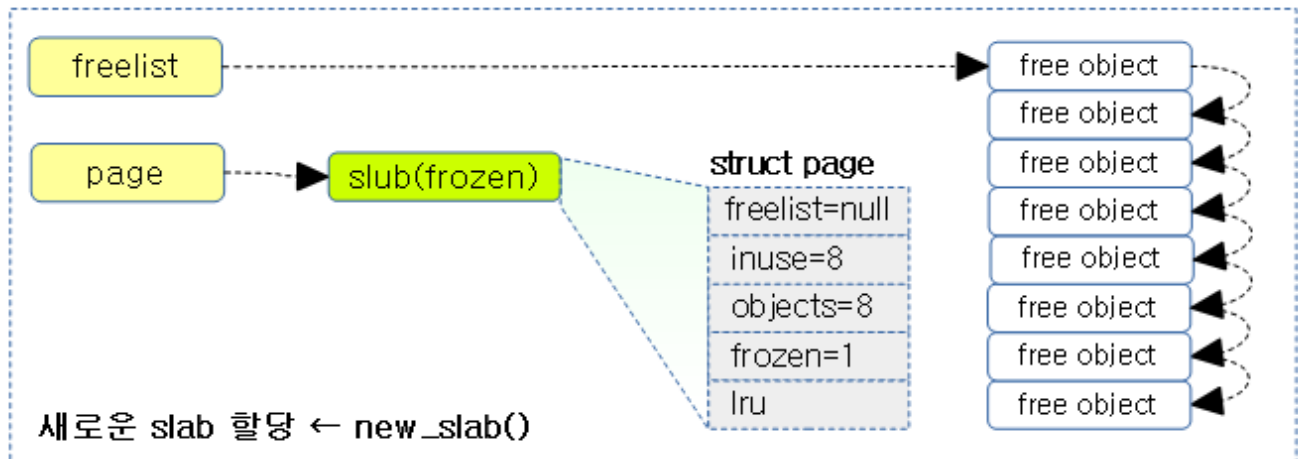
## Inuse counter on the slub page

The inuse counter on the slub page indicates the number of objects in use, but also includes the free objects in the c->freelist as the number of objects in use. The changes to the inuse counter are as follows:

- Initial value
  - After the slab page is first assigned, all free objects are linked to C->freelist, so the number of page->objects is assigned to the page->inuse value to mark them as all in use.
- At the time of assignment
  - When you assign an object, there is no change to the page->inuse counter.
- On Release
  - When freed from this CPU, the page->inuse counter remains unchanged.
  - When disabling from the remote CPU, the page->inuse counter is decremented.
- On the Move
  - When moving from c->page->freelist to c->freelist, the page->inuse value is reflected as the page->objects value, and all objects are in use.

**unfrozen state**

Free objects in slab pages managed by n->partial are not included in the number of objects in use.

The following three figures show that the Page->inuse counter is decremented only when deallocating from the remote CPU.



(http://jake.dothome.co.kr/wp-content/uploads/2016/06/frozen-slub-page-3.png)



(http://jake.dothome.co.kr/wp-content/uploads/2016/06/frozen-slub-page-4.png)



(http://jake.dothome.co.kr/wp-content/uploads/2016/06/frozen-slub-page-5a.png)

# Assign a slab cache object

The following figure outlines the order in which the allocation request for a slub object is performed by function. If there are no objects to allocate by the last step, the Buddy system allocates a page to construct a slub and initializes the objects inside it.



(http://jake.dothome.co.kr/wp-content/uploads/2016/06/kmem_cache_alloc-1a.png)
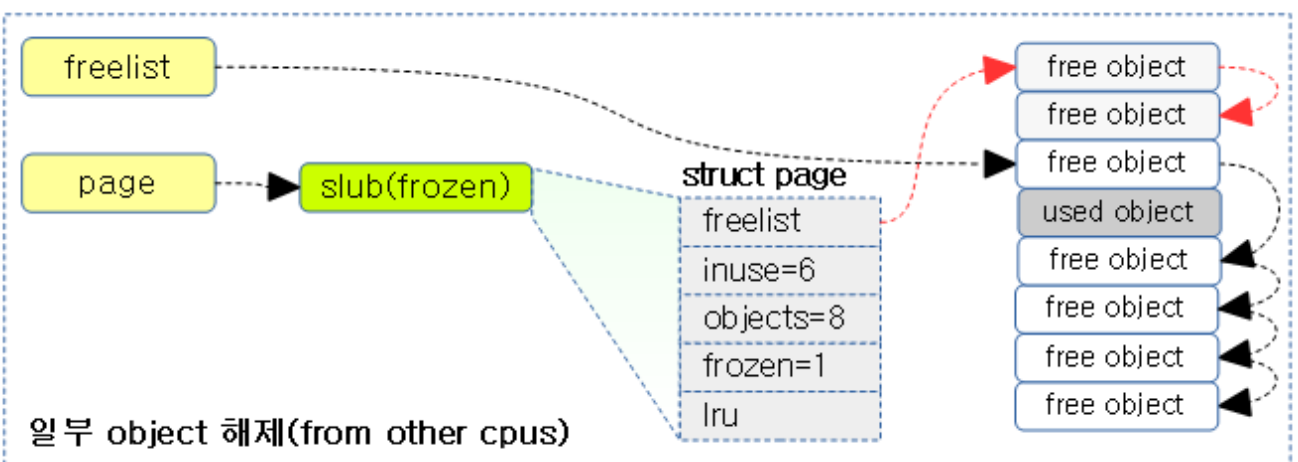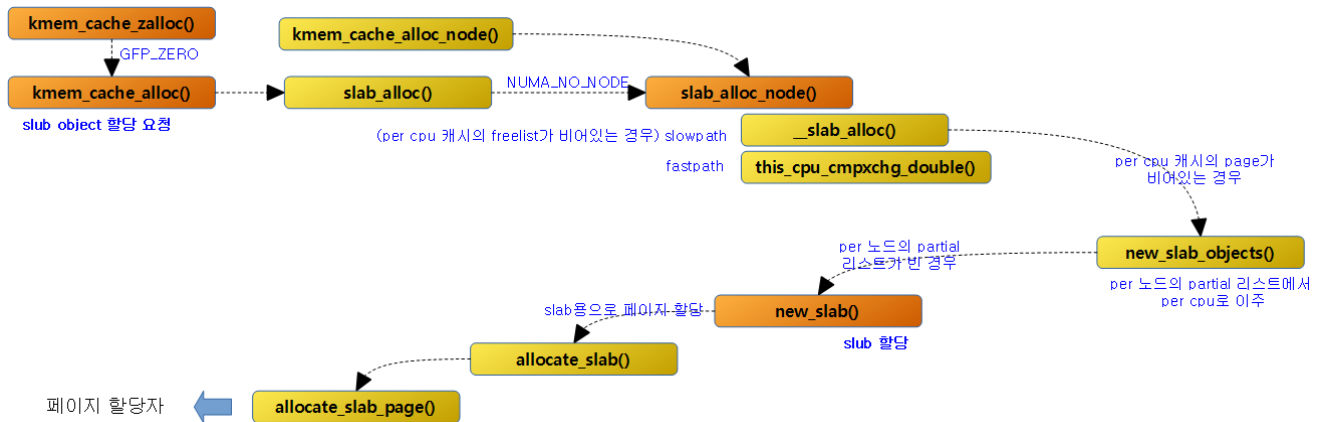
## kmem_cache_zalloc()

include/linux/slab.h

```
1  /*
2   * Shortcuts
3   */
```

```
1  static inline void *kmem_cache_zalloc(struct kmem_cache *k, gfp_t flags)
2  {
3          return kmem_cache_alloc(k, flags | __GFP_ZERO);
4  }
```

It is allocated a slub object initialized to zero from the specified cache.

## kmem_cache_alloc()

mm/slub.c

```
01  void *kmem_cache_alloc(struct kmem_cache *s, gfp_t gfpflags)
02  {
03          void *ret = slab_alloc(s, gfpflags, _RET_IP_);
04
05          trace_kmem_cache_alloc(_RET_IP_, ret, s->object_size,
06                                 s->size, gfpflags);
07
08          return ret;
09  }
10  EXPORT_SYMBOL(kmem_cache_alloc);
```

It is allocated a slub object from the specified cache.

- trace 문자열: "call_site=%lx ptr=%p bytes_req=%zu bytes_alloc=%zu gfp_flags=%s"

## kmem_cache_alloc_node()

mm/slub.c

```
01  void *kmem_cache_alloc_node(struct kmem_cache *s, gfp_t gfpflags, int no
    de)
02  {
03          void *ret = slab_alloc_node(s, gfpflags, node, _RET_IP_);
04
05          trace_kmem_cache_alloc_node(_RET_IP_, ret,
06                                      s->object_size, s->size, gfpflags, n
    ode);
07
08          return ret;
09  }
10  EXPORT_SYMBOL(kmem_cache_alloc_node);
```

GFP flag to allocate a slub object from the specified node.

## slab_alloc()

mm/slub.c

```
1  static __always_inline void *slab_alloc(struct kmem_cache *s,
2                  gfp_t gfpflags, unsigned long addr)
3  {
4          return slab_alloc_node(s, gfpflags, NUMA_NO_NODE, addr);
5  }
```

One slub object is assigned to the entire node from the specified cache.

# Assign a slab object – Fastpath

## slab_alloc_node() – slub

mm/slub.c

```
01  /*
02   * Inlined fastpath so that allocation functions (kmalloc, kmem_cache_al
    loc)
03   * have the fastpath folded into their functions. So no function call
04   * overhead for requests that can be satisfied on the fastpath.
05   *
06   * The fastpath works by first checking if the lockless freelist can be
    used.
07   * If not then __slab_alloc is called for slow processing.
08   *
09   * Otherwise we can simply pick the next object from the lockless free l
    ist.
10   */
```

```
01  static __always_inline void *slab_alloc_node(struct kmem_cache *s,
02                  gfp_t gfpflags, int node, unsigned long addr)
03  {
04          void **object;
05          struct kmem_cache_cpu *c;
06          struct page *page;
07          unsigned long tid;
08
09          s = slab_pre_alloc_hook(s, gfpflags);
10          if (!s)
11                  return NULL;
```

```
12  redo:
13          /*
14           * Must read kmem_cache cpu data via this cpu ptr. Preemption is
15           * enabled. We may switch back and forth between cpus while
16           * reading from one cpu area. That does not matter as long
17           * as we end up on the original cpu again when doing the cmpxch
    g.
18           *
19           * We should guarantee that tid and kmem_cache are retrieved on
20           * the same cpu. It could be different if CONFIG_PREEMPT so we n
    eed
21           * to check if it is matched or not.
22           */
23          do {
24                  tid = this_cpu_read(s->cpu_slab->tid);
25                  c = raw_cpu_ptr(s->cpu_slab);
26          } while (IS_ENABLED(CONFIG_PREEMPT) &&
27                   unlikely(tid != READ_ONCE(c->tid)));
28
29          /*
30           * Irqless object alloc/free algorithm used here depends on sequ
    ence
31           * of fetching cpu_slab's data. tid should be fetched before any
    thing
32           * on c to guarantee that object and page associated with previo
    us tid
33           * won't be used with current tid. If we fetch tid first, object
    and
34           * page could be one associated with next tid and our alloc/free
35           * request will be failed. In this case, we will retry. So, no p
    roblem.
36           */
37          barrier();
38
39          /*
40           * The transaction ids are globally unique per cpu and per opera
    tion on
41           * a per cpu queue. Thus they can be guarantee that the cmpxchg_
    double
42           * occurs on the right processor and that there was no operation
    on the
43           * linked list in between.
44           */
45
46          object = c->freelist;
47          page = c->page;
48          if (unlikely(!object || !node_match(page, node))) {
49                  object = __slab_alloc(s, gfpflags, node, addr, c);
50                  stat(s, ALLOC_SLOWPATH);
51          } else {
52                  void *next_object = get_freepointer_safe(s, object);
53
54                  /*
55                   * The cmpxchg will only match if there was no additiona
    l
56                   * operation and if we are on the right processor.
57                   *
58                   * The cmpxchg does the following atomically (without lo
    ck
59                   * semantics!)
60                   * 1. Relocate first pointer to the current per cpu are
    a.
61                   * 2. Verify that tid and freelist have not been changed
62                   * 3. If they were not changed replace tid and freelist
63                   *
64                   * Since this is without lock semantics the protection i
    s only
```

```
65                          * against code executing on this cpu *not* from access
   by
66                          * other cpus.
67                          */
68                      if (unlikely(!this_cpu_cmpxchg_double(
69                                      s->cpu_slab->freelist, s->cpu_slab->tid,
70                                      object, tid,
71                                      next_object, next_tid(tid)))) {
72
73                              note_cmpxchg_failure("slab_alloc", s, tid);
74                              goto redo;
75                      }
76                      prefetch_freepointer(s, next_object);
77                      stat(s, ALLOC_FASTPATH);
78              }
79
80          if (unlikely(gfpflags & __GFP_ZERO) && object)
81              memset(object, 0, s->object_size);
82
83          slab_post_alloc_hook(s, gfpflags, 1, &object);
84
85          return object;
86 }
```

It is allocated a slub object from the specified node in the specified cache.

- In line 9~11 of the code, we perform the pre-preparations for the slab object assignment, and if there is a problem, we give up the assignment. If the memcg cache is in effect, the memcg cache is obtained.
  - If there is a problem with the allocation of slab memory due to fault injection, the allocation is abandoned.
- On line 12 of the code, the redo: label is. For fastpath purposes, the transaction ID for that CPU has changed from the per CPU cache, and if it fails, it will repeat again.
- In lines 23~27 of the code, the TID and Per-CPU caches are read atomically. At this point, preemption can be done at any time, so it is possible to switch tasks and come back during execution. Therefore, in this routine, the verification process is performed behind the scenes via atomic to ensure that TID and cache are obtained from the same CPU.
- In line 37 of code, in order for slab's allocation/release algorithm to work without interrupt masks, it relies on the order in which cpu_slab data is read. In order to read tid before object and page, I used a compiler barrier to clearly distinguish the order of operations so that the compiler would not have to do optimization.
- In line 46~47 of the code, we get an object from the freelist of the per cpu cache, and we also get a page.
- In line 48~50 of the code, if there is a low probability that there is no slab object to allocate, or if the node of the slab page known from the current per CPU cache is different from the specified node, the object will be assigned via the slowpath operation. At this point, the ALLOC_SLOWPATH stat increases.
- If the object is successfully fetched from lines 51~52 of the code, the FP(Free Pointer) value, i.e., the address of the next object, is known.
- In lines 68~75 of the code, try to operate the fastpath operation of the slab object assignment. Atomically, the freelist points to the next slab object, and the tid is replaced with the next tid. If there is a small probability that it fails due to a change in the CPU or a change in the transaction ID, it will revert back to the redo label and retry until the slab object allocation succeeds.

- If a task is preemption and resumed on a different node's CPU, the transaction ID may be changed by the CPU or by a different remote CPU. In this case, give up using the object obtained from the freelist and try again.
  - 참고: Per-cpu (atomic operations) (http://jake.dothome.co.kr/per-cpu-atomic) | 문c
- In line 76~77 of the code, the assignment of the slab object was successful with the fastpath operation. Prefetch the cacheline for the next free object, and increment the ALLOC_FASTPATH stat.
- If you use the GFP_ZERO flag with a small probability in line 80~81 of the code, clear the object to 0 by the size of the assigned object.
- In line 83 of code, the assigned completed state is notified to debug (kmemcheck, kmemleak, kasan) routines, etc., and to memcg.
- Returns the slab object assigned in line 85 of code.

## Lock-less per-cpu list manipulation

In general, locks are used for list manipulation, but for fast list manipulation without using locks, it is often implemented with RCU or Per-CPU. Since the allocation and decommissioning of slab objec is frequently used in kernel code, fast performance is also required, so the method using per-cpu was chosen among the solutions that support lock-less implementation. If you implement using per-CPU, you typically dispose of interrupts in order to block preemption and interrupt intervention. However, the allocation and deallocation of slab objects is implemented with preemption and interrupts allowed for faster performance. Using this technique, it is necessary to ensure that there is no problem with list manipulation in order to handle the situation where the slab object is preempted during the allocation and deallocation process. Therefore, it is possible to safely handle list manipulation using atomic, which can be exchanged after comparison (cmpxchg), but if only the pointer of an object is compared, it may not be possible to handle it safely in the event of a special situation such as an ABA problem. This has led to the need for additional transaction IDs and features such as cpmxchg_double to compare and process two values at once, which are supported by the latest x86_64 and arm64-bit architectures.

### this_cpu_cmpxchg()

- this_cpu_cmpxchg(v, old, new)
- Compare one v value with the old value and change it to a new value only if it is the same.

### this_cpu_cmpxcgh_double()

- this_cpu_cmpxcgh_double(v1, v2, old1, old2, new1, new2)
- Compare the two v1 and v2 values with the old1 and old2 values and change them to new1 and new2 values for the year, provided they are both the same.

### Transaction ID

- Transaction IDs have different starting values on each CPU so that they don't overlap with each other, and use larger values to avoid intersecting incremental values.
- If you are using per-cpu with preemption enabled, it is safe to use the transaction id to check whether the cpu has changed, as the value of the currently accessed per-cpu may change due to changes in the cpu.

- If you try to change it with only the object address and not the transaction ID, you may have an ABA problem.
  - Consultation:
    - Multiprocessing and Multicore (2013) | Timothy Roscoe – Download PDF (https://www.systems.ethz.ch/sites/default/files/file/COURSES/2013_FALL_COURSES/ 2013_Fall_SPCA/lectures/19-Multiprocessing-1up.pdf)
    - Avoid the ABA problem when using lock-free algorithms (https://wiki.sei.cmu.edu/confluence/display/c/CON09- C.+Avoid+the+ABA+problem+when+using+lock-free+algorithms) | Carnegie Mellon University

The following figure shows how a slab object allocation request quickly allocates the first object from the freelist in the per CPU cache.



(http://jake.dothome.co.kr/wp-content/uploads/2016/06/slab_alloc_node-1f.png)

## slab_pre_alloc_hook()

mm/slub.c

```
01  static inline struct kmem_cache *slab_pre_alloc_hook(struct kmem_cache *
    s,
02                                                     gfp_t flags)
03  {
04          flags &= gfp_allowed_mask;
05
06          fs_reclaim_acquire(flags);
07          fs_reclaim_release(flags);
08
09          might_sleep_if(gfpflags_allow_blocking(flags));
10
11          if (should_failslab(s, flags))
12                  return NULL;
13
14          if (memcg_kmem_enabled() &&
15              ((flags & __GFP_ACCOUNT) || (s->flags & SLAB_ACCOUNT)))
16                  return memcg_kmem_get_cache(s);
17
18          return s;
19  }
```

Perform the prerequisites for the slab object assignment, and if there is a problem, give up the assignment. It also returns the memcg cache if it is in effect.

- In line 4 of code, make sure that only the allowed GFP flags are used on early boot.
- If you're doing a direct-reclaim on line 9 of code, trigger the preemption point.

- Sleep the current task and hand over the CPU to another task to handle requests with a higher priority than the current task.
- If you use the CONFIG_FAILSLAB kernel option in lines 11~12 of code, it will return null to give up the allocation if you are forced to commit a failure by the fault-injection function.
- In the case of a task with memcg applied in line 14~16 of code, find the per memcg cache and return it.

## memcg_kmem_get_cache()

mm/memcontrol.c

```
01  /**
02   * memcg_kmem_get_cache: select the correct per-memcg cache for allocati
     on
03   * @cachep: the original global kmem cache
04   *
05   * Return the kmem_cache we're supposed to use for a slab allocation.
06   * We try to use the current memcg's version of the cache.
07   *
08   * If the cache does not exist yet, if we are the first user of it, we
09   * create it asynchronously in a workqueue and let the current allocatio
     n
10   * go through with the original cache.
11   *
12   * This function takes a reference to the cache it returns to assure it
13   * won't get destroyed while we are working with it. Once the caller is
14   * done with it, memcg_kmem_put_cache() must be called to release the
15   * reference.
16   */
```

```
01  struct kmem_cache *memcg_kmem_get_cache(struct kmem_cache *cachep)
02  {
03          struct mem_cgroup *memcg;
04          struct kmem_cache *memcg_cachep;
05          int kmemcg_id;
06
07          VM_BUG_ON(!is_root_cache(cachep));
08
09          if (memcg_kmem_bypass())
10                  return cachep;
11
12          memcg = get_mem_cgroup_from_current();
13          kmemcg_id = READ_ONCE(memcg->kmemcg_id);
14          if (kmemcg_id < 0)
15                  goto out;
16
17          memcg_cachep = cache_from_memcg_idx(cachep, kmemcg_id);
18          if (likely(memcg_cachep))
19                  return memcg_cachep;
20
21          /*
22           * If we are in a safe context (can wait, and not in interrupt
23           * context), we could be be predictable and return right away.
24           * This would guarantee that the allocation being performed
25           * already belongs in the new cache.
26           *
27           * However, there are some clashes that can arrive from locking.
28           * For instance, because we acquire the slab_mutex while doing
29           * memcg_create_kmem_cache, this means no further allocation
30           * could happen with the slab_mutex held. So it's better to
31           * defer everything.
32           */
33          memcg_schedule_kmem_cache_create(memcg, cachep);
```

```
34   out:
35           css_put(&memcg->css);
36           return cachep;
37   }
```

If the current user task has memcg enabled, select the per memcg cache.

- If interrupted, current->mm is null, or if you are running as a kernel thread, do not change cachep.

## slab_post_alloc_hook()

mm/slab.h

```
01   static inline void slab_post_alloc_hook(struct kmem_cache *s, gfp_t flag
     s,
02                                            size_t size, void **p)
03   {
04           size_t i;
05
06           flags &= gfp_allowed_mask;
07           for (i = 0; i < size; i++) {
08                   p[i] = kasan_slab_alloc(s, p[i], flags);
09                   /* As p[i] might get tagged, call kmemleak hook after KA
     SAN. */
10                   kmemleak_alloc_recursive(p[i], s->object_size, 1,
11                                            s->flags, flags);
12           }
13
14           if (memcg_kmem_enabled())
15                   memcg_kmem_put_cache(s);
16   }
```

Perform a routine for post-performance in the slab object assignment. Perform debugging routines such as KASAN, kmemleak, etc.

## node_match()

mm/slub.c

```
1   /*
2    * Check if the objects in a per cpu structure fit numa
3    * locality expectations.
4    */
```

```
1   static inline int node_match(struct page *page, int node)
2   {
3   #ifdef CONFIG_NUMA
4           if (!page || (node != NUMA_NO_NODE && page_to_nid(page) != nod
    e))
5                   return 0;
6   #endif
7           return 1;
8   }
```

In a UMA system, it always returns true, and in a NUMA system, it returns false if the node in page(slub) specified as an argument is different from the node specified as an argument.

- Using the preempt kernel of the NUMA system, the slub allocation routine should be preempted when attempting to allocate from the requested node, which will receive an object from the page

(slub) of the current per cpu cache, and compare it to the node belonging to that page(slub) and return false in any other case so that it can be determined that the current cpu cache should not be used as is.

## Free Pointer (FP)

### prefetch_freepointer()

mm/slub.c

```
1  static void prefetch_freepointer(const struct kmem_cache *s, void *objec
   t)
2  {
3          prefetch(object + s->offset);
4  }
```

Load the FP(Freelist Pointer) value of the requested slab object into the cache line in advance.

### get_freepointer_safe()

mm/slub.c

```
01  static inline void *get_freepointer_safe(struct kmem_cache *s, void *obj
    ect)
02  {
03          unsigned long freepointer_addr;
04          void *p;
05
06          if (!debug_pagealloc_enabled())
07                  return get_freepointer(s, object);
08
09          freepointer_addr = (unsigned long)object + s->offset;
10          probe_kernel_read(&p, (void **)freepointer_addr, sizeof(p));
11          return freelist_ptr(s, p, freepointer_addr);
12  }
```

It reads the FP(Freelist Pointer) value of the requested slab object and converts it to a pointer address.

### get_freepointer()

mm/slub.c

```
1  static inline void *get_freepointer(struct kmem_cache *s, void *object)
2  {
3          return freelist_dereference(s, object + s->offset);
4  }
```

Reads the FP (Freelist Pointer) value of the specified object and returns it to its address.

### freelist_dereference()

mm/slub.c

```
1  /* Returns the freelist pointer recorded at location ptr_addr. */
2  static inline void *freelist_dereference(const struct kmem_cache *s,
```

```
  3                                    void *ptr_addr)
  4 {
  5         return freelist_ptr(s, (void *)*(unsigned long *)(ptr_addr),
  6                              (unsigned long)ptr_addr);
  7 }
```

Converts the @ptr_addr value, which is the FP (Freelist Pointer) value, into a pointer address and returns it.

### freelist_ptr()

mm/slub.c

```
  1 /*
  2  * Returns freelist pointer (ptr). With hardening, this is obfuscated
  3  * with an XOR of the address where the pointer is held and a per-cache
  4  * random number.
  5  */
```

```
 01 static inline void *freelist_ptr(const struct kmem_cache *s, void *ptr,
 02                                  unsigned long ptr_addr)
 03 {
 04 #ifdef CONFIG_SLAB_FREELIST_HARDENED
 05         /*
 06          * When CONFIG_KASAN_SW_TAGS is enabled, ptr_addr might be tagged.
 07          * Normally, this doesn't cause any issues, as both set_freepointer()
 08          * and get_freepointer() are called with a pointer with the same tag.
 09          * However, there are some issues with CONFIG_SLUB_DEBUG code. For
 10          * example, when __free_slub() iterates over objects in a cache, it
 11          * passes untagged pointers to check_object(). check_object() in turns
 12          * calls get_freepointer() with an untagged pointer, which causes the
 13          * freepointer to be restored incorrectly.
 14          */
 15         return (void *)((unsigned long)ptr ^ s->random ^
 16                         (unsigned long)kasan_reset_tag((void *)ptr_addr));
 17 #else
 18         return ptr;
 19 #endif
 20 }
```

Converts the @ptr_addr value, which is the FP (Freelist Pointer) value, into a pointer address and returns it.

- If you use the CONFIG_SLAB_FREELIST_HARDENED kernel options, the FP value is encoded along with a random value so that it is not recognizable.

## Assign a Slap Object – Slowpath

### __slab_alloc()

mm/slub.c

```
 1  /*
 2   * Another one that disabled interrupt and compensates for possible
 3   * cpu changes by refetching the per cpu area pointer.
 4   */

01  static void *__slab_alloc(struct kmem_cache *s, gfp_t gfpflags, int nod
    e,
02                             unsigned long addr, struct kmem_cache_cpu *c)
03  {
04          void *p;
05          unsigned long flags;
06
07          local_irq_save(flags);
08  #ifdef CONFIG_PREEMPT
09          /*
10           * We may have been preempted and rescheduled on a different
11           * cpu before disabling interrupts. Need to reload cpu area
12           * pointer.
13           */
14          c = this_cpu_ptr(s->cpu_slab);
15  #endif
16
17          p = ___slab_alloc(s, gfpflags, node, addr, c);
18          local_irq_restore(flags);
19          return p;
20  }
```

With interrupts on the local CPU blocked, perform a slowpath routine for slab object assignment.


## \_\_\_slab_alloc()

mm/slub.c

```
01  /*
02   * Slow path. The lockless freelist is empty or we need to perform
03   * debugging duties.
04   *
05   * Processing is still very fast if new objects have been freed to the
06   * regular freelist. In that case we simply take over the regular freeli
    st
07   * as the lockless freelist and zap the regular freelist.
08   *
09   * If that is not working then we fall back to the partial lists. We tak
    e the
10   * first element of the freelist as the object to allocate now and move
    the
11   * rest of the freelist to the lockless freelist.
12   *
13   * And if we were unable to get a new slab from the partial slab lists t
    hen
14   * we need to allocate a new slab. This is the slowest path since it inv
    olves
15   * a call to the page allocator and the setup of a new slab.
16   *
17   * Version of __slab_alloc to use when we know that interrupts are
18   * already disabled (which is the case for bulk allocation).
19   */

01  static void *___slab_alloc(struct kmem_cache *s, gfp_t gfpflags, int nod
    e,
02                             unsigned long addr, struct kmem_cache_cpu *c)
03  {
04          void *freelist;
05          struct page *page;
06
```

```
07          page = c->page;
08          if (!page)
09                  goto new_slab;
10  redo:
11
12          if (unlikely(!node_match(page, node))) {
13                  int searchnode = node;
14
15                  if (node != NUMA_NO_NODE && !node_present_pages(node))
16                          searchnode = node_to_mem_node(node);
17
18                  if (unlikely(!node_match(page, searchnode))) {
19                          stat(s, ALLOC_NODE_MISMATCH);
20                          deactivate_slab(s, page, c->freelist, c);
21                          goto new_slab;
22                  }
23          }
24
25          /*
26           * By rights, we should be searching for a slab page that was
27           * PFMEMALLOC but right now, we are losing the pfmemalloc
28           * information when the page leaves the per-cpu allocator
29           */
30          if (unlikely(!pfmemalloc_match(page, gfpflags))) {
31                  deactivate_slab(s, page, c->freelist, c);
32                  goto new_slab;
33          }
34
35          /* must check again c->freelist in case of cpu migration or IRQ
    */
36          freelist = c->freelist;
37          if (freelist)
38                  goto load_freelist;
39
40          freelist = get_freelist(s, page);
41
42          if (!freelist) {
43                  c->page = NULL;
44                  stat(s, DEACTIVATE_BYPASS);
45                  goto new_slab;
46          }
47
48          stat(s, ALLOC_REFILL);
49
50  load_freelist:
51          /*
52           * freelist is pointing to the list of objects to be used.
53           * page is pointing to the page from which the objects are obtai
    ned.
54           * That page must be frozen for per cpu allocations to work.
55           */
56          VM_BUG_ON(!c->page->frozen);
57          c->freelist = get_freepointer(s, freelist);
58          c->tid = next_tid(c->tid);
59          return freelist;
60
61  new_slab:
62
63          if (slub_percpu_partial(c)) {
64                  page = c->page = slub_percpu_partial(c);
65                  slub_set_percpu_partial(c, page);
66                  stat(s, CPU_PARTIAL_ALLOC);
67                  goto redo;
68          }
69
70          freelist = new_slab_objects(s, gfpflags, node, &c);
71
72          if (unlikely(!freelist)) {
```

```
73                    slab_out_of_memory(s, gfpflags, node);
74                    return NULL;
75              }
76
77              page = c->page;
78              if (likely(!kmem_cache_debug(s) && pfmemalloc_match(page, gfpfla
    gs)))
79                    goto load_freelist;
80
81              /* Only entered in the debug case */
82              if (kmem_cache_debug(s) &&
83                          !alloc_debug_processing(s, page, freelist, add
    r))
84                    goto new_slab;  /* Slab failed checks. Next slab needed
    */
85
86              deactivate_slab(s, page, get_freepointer(s, freelist), c);
87              return freelist;
88    }
```
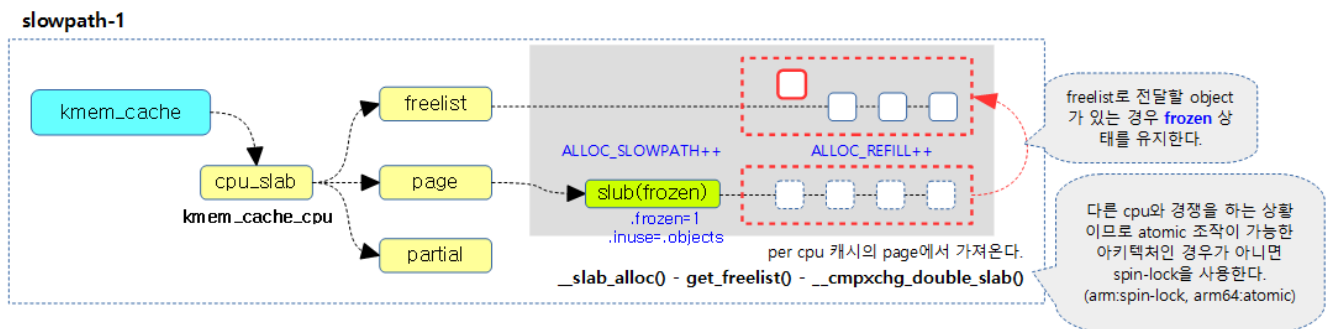
Perform a step-by-step slowpath routine for slab object assignment. s->page->freelist –> s->partial –>
n->partial steps and if that doesn't work, the last step is to allocate slab pages directly from the page
allocator.

- If c->page is not specified in code lines 7~9, go to the new_slab to fetch a new slab page from the
  c->partial list.
- On line 10 of the code, the redo: label is. It is a label that is moved after charging a slab page
  from a c->partial or n->partial list to return to this label and retry the object assignment.
- In line 12~23 of the code, there is a low probability that the node in c->page and the node
  requested as an argument are different. It can also be a memoryless node configuration, so if
  the node in c->page is different from the adjacent node again, increment the
  ALLOC_NODE_MISMATCH counter, deactivate the slab page linked to s->page, assign null to page
  and freelist, and then go to the new_slab to get the new slab page from the c->partial list.
    - Normally, if a task is preemptioned and resched, and the task resumes from that location,
      it will be rescheduled to the original CPU or CPU belonging to that node, if possible, but if
      the CPU of that node is busy, it may be rescheduled to the CPU of another node. This is the
      case when the node on the current CPU may be different from the node that was originally
      requested.
- In line 30~33 of code, if the current request is not an emergency and an emergency slub page is
  selected, deactivate the slab page linked to s->page, assign null to page and freelist, and move to
  the new_slab to fetch a new slub page from the c->partial list.
    - If you are using the swap for NBD driver, use the ALLOC_NO_WATERMARKS flag when
      requesting a slab object assignment to allow TCP socket senders and receivers to take
      advantage of the emergency slub page as well.
- In lines 36~38 of the code, once again check that c->freelist is empty. If the object returned just
  before the interrupt is left in the freelist, it moves to the load_freelist label to receive the slab
  object from c->freelist.
- On line 40, move c->page->freelist to c->freelist. (SlowPath Step 1)
- In code lines 42~46, if there are no free slab objects in the c->freelist, the slab page is in use, so
  we assign null to c->page to quietly stop the slab page from being managed by the slab allocator,
  and increment the DEACTIVATE_BYPASS counter. Then go to the new_slab label to get a new slap
  page assigned. (This is the case when the first step of the slowpath failed.)

- In line 48 of the code, c->freelist is repopulated, so increment the ALLOC_REFILL counter. (This is the case if the Slowpath step 1 succeeds.)
- In code lines 50~59, load_freelist: Label. c->freelist points to the next free object, increments the transaction id, and returns the assigned slab object.
- In code line 61, new_slab: Label.
- In code lines 63~68, if a slab page exists in c->partial, only one slab page is moved to c->page. Then increment the CPU_PARTIAL_ALLOC counter and go to the redo: label to retry the allocation of the slab object. (SlowPath Step 2)
- On line 70 of the code, new_slab_objects you want to be allocated a new slab page from an n->partial list or buddy allocator, call the function to connect to c->page and c->freelist and get a new free object assigned. (Slowpath 3~4 steps)
- If a new slab object is not allocated in line 72~75, it will return null after OOM (Out Of Memory).
- In line 77~79 of the code, if all of the following situations are met with a high probability, go to the load_freelist label and retry the free object assignment.
  - When you don't need to debug Slub
  - The current request is not an emergency, but the emergency slub page is selected.
- If a debug request SLAB_CONSISTENCY_CHECKS in line 82~84 of code finds a problem in the consistency check, it moves to the new_slab label to be assigned a new slab page.
- In line 86 of code, deactivate the slab page linked to s->page and assign null to page and freelist.
- Returns the slab object assigned in line 87 of code.

The following diagram shows the SlowPath step 1 process of moving from c->page->freelist to c->freelist>



(http://jake.dothome.co.kr/wp-content/uploads/2016/06/slab_alloc_node-2c.png)

The following diagram shows the two-step Slowpath process of moving to c->partial – > c->page & c->freelist.

- Since all the objects on the slab page that were previously managed have been allocated, they are not directly managed by the slab cache allocator.

(http://jake.dothome.co.kr/wp-content/uploads/2016/06/slab_alloc_node-3d.png)

## Slowpath Step 1 Processing

### get_freelist()

mm/slub.c

```
01   /*
02    * Check the page->freelist of a page and either transfer the freelist to the
03    * per cpu freelist or deactivate the page.
04    *
05    * The page is still frozen if the return value is not NULL.
06    *
07    * If this function returns NULL then the page has been unfrozen.
08    *
09    * This function must be called with interrupt disabled.
10    */
```

```
01   static inline void *get_freelist(struct kmem_cache *s, struct page *page)
02   {
03           struct page new;
04           unsigned long counters;
05           void *freelist;
06
07           do {
08                   freelist = page->freelist;
09                   counters = page->counters;
10
11                   new.counters = counters;
12                   VM_BUG_ON(!new.frozen);
13
14                   new.inuse = page->objects;
15                   new.frozen = freelist != NULL;
16
17           } while (!__cmpxchg_double_slab(s, page,
18                           freelist, counters,
```

```
19              NULL, new.counters,
20              "get_freelist"));
21
22      return freelist;
23 }
```

Returns the freelist of the page in the per CPU cache, and puts null into the freelist of the slub page, and changes it to All in use. (SlowPath Step 1)

- c->page->freelist – > When all free objects are moved to C->freelist, the value of S->Objects is assigned to the inuse containing the number of objects in use and changed to all of them in use.
- Assign null to page->freelist, and new.counters to page->counter, and repeat if that fails

## Slowpath 3 & Step 4 Processing

### new_slab_objects()

mm/slub.c

```
01 static inline void *new_slab_objects(struct kmem_cache *s, gfp_t flags,
02                      int node, struct kmem_cache_cpu **pc)
03 {
04      void *freelist;
05      struct kmem_cache_cpu *c = *pc;
06      struct page *page;
07
08      WARN_ON_ONCE(s->ctor && (flags & __GFP_ZERO));
09      freelist = get_partial(s, flags, node, c);
10
11      if (freelist)
12              return freelist;
13
14      page = new_slab(s, flags, node);
15      if (page) {
16              c = raw_cpu_ptr(s->cpu_slab);
17              if (c->page)
18                      flush_slab(s, c);
19
20              /*
21               * No other reference to the page yet so we can
22               * muck around with it freely without cmpxchg
23               */
24              freelist = page->freelist;
25              page->freelist = NULL;
26
27              stat(s, ALLOC_SLAB);
28              c->page = page;
29              *pc = c;
30      } else
31              freelist = NULL;
32
33      return freelist;
34 }
```

Move a few slab pages from the n->partial list to the c->freelist, c->page, and c->partial lists. If there is no slab page to move, a new slab page is assigned by the buddy system and assigned to c->freelist and c->page. If even page allocation fails, it returns null. (Slowpath 3~4 steps)

- In line 10~13 of the code, move a few slab pages from the partial list of local and remote nodes to the c->freelist, c->page, and c->partial lists, and return if there is a slab object to be assigned normally. (Slowpath Step 3)
- In line 15~30 of the code, a new page is assigned from the buddy system. It then connects to c->page, inserts nulls into c->page->freelist, and increments the ALLOC_SLAB counter. Since the CPU can be changed, the per-CPU cache is substituted in the output argument @pc. It then returns a freelist. (Slowpath Step 4)
- In line 31~32 of the code, the buddy system failed to allocate due to insufficient memory. NULL.

The following diagram shows the three-step slowpath process of moving to n->partial – > c->page & c->freelist & c->partial.

- If even this node is insufficient, it searches for the remote node and performs it.



(http://jake.dothome.co.kr/wp-content/uploads/2016/06/slab_alloc_node-4f.png)

The following diagram shows the SlowPath 4-step process of moving a slap page from the Buddy system to C->Page & C->Freelist if there is a shortage of slab pages.

새로운 페이지를 할당받아서 per cpu 캐시로 이동

(http://jake.dothome.co.kr/wp-content/uploads/2016/06/slab_alloc_node-5b.png)

## SlowPath Step 3 Processing Call

### get_partial()

mm/slub.c

```
 1  /*
 2   * Get a partial page, lock it and return it.
 3   */

01  static void *get_partial(struct kmem_cache *s, gfp_t flags, int node,
02                  struct kmem_cache_cpu *c)
03  {
04          void *object;
05          int searchnode = node;
06
07          if (node == NUMA_NO_NODE)
08                  searchnode = numa_mem_id();
09          else if (!node_present_pages(node))
10                  searchnode = node_to_mem_node(node);
11
12          object = get_partial_node(s, get_node(s, searchnode), c, flags);
13          if (object || node != NUMA_NO_NODE)
14                  return object;
15
16          return get_any_partial(s, flags, c);
17  }
```

Move the regular slab pages of the n->partial list to the c->freelist, c->page, and c->partial lists. If the partial list of specified nodes is empty, try on the remote nodes. (Slowpath Step 3)

- Migrate the slab pages from the partial list of nodes specified in code lines 7~14 (adjacent nodes in the case of memoryless nodes) to page and partial in the per CPU cache. If there is a free object to be allocated normally, it returns the object.
- In line 16 of code, we perform the same processing as above with remote nodes.

## Slowpath Step 3 – Processing on the Request Node

### get_partial_node()

mm/slub.c

```
1   /*
2    * Try to allocate a partial slab from a specific node.
3    */

01  static void *get_partial_node(struct kmem_cache *s, struct kmem_cache_no
    de *n,
02                                      struct kmem_cache_cpu *c, gfp_t flags)
03  {
04          struct page *page, *page2;
05          void *object = NULL;
06          int available = 0;
07          int objects;
08
09          /*
10           * Racy check. If we mistakenly see no partial slabs then we
11           * just allocate an empty slab. If we mistakenly try to get a
12           * partial slab and there is none available then get_partials()
13           * will return NULL.
14           */
15          if (!n || !n->nr_partial)
16                  return NULL;
17
18          spin_lock(&n->list_lock);
19          list_for_each_entry_safe(page, page2, &n->partial, lru) {
20                  void *t;
21
22                  if (!pfmemalloc_match(page, flags))
23                          continue;
24
25                  t = acquire_slab(s, n, page, object == NULL, &objects);
26                  if (!t)
27                          break;
28
29                  available += objects;
30                  if (!object) {
31                          c->page = page;
32                          stat(s, ALLOC_FROM_PARTIAL);
33                          object = t;
34                  } else {
35                          put_cpu_partial(s, page, 0);
36                          stat(s, CPU_PARTIAL_NODE);
37                  }
38                  if (!kmem_cache_has_cpu_partial(s)
39                          || available > s->cpu_partial / 2)
40                          break;
41
42          }
43          spin_unlock(&n->list_lock);
44          return object;
45  }
```

It retrieves a certain amount of slap pages from a partial list of per nodes and migrates them to the per CPU cache. Returns null with no slab pages to migrate. (Slowpath Step 3 – Local Node)

- In line 15~16 of the code, if the number of slab pages in the per node reaches zero, it returns null.
  - Since there is no node lock, it can be zero in a node lock competition.

- With the node lock obtained in line 18~23 of code, it traverses the slab pages in the partial list of requested nodes, and pfmemalloc skips the unmatched slab pages.
- If the first slap page is fetched in lines 25~33, move it to page in the per CPU cache, and increment the ALLOC_FROM_PARTIAL counter.
- Starting from the second slap page in lines 34~37 of code, move it to the partial list of per cpu cache, and increment the CPU_PARTIAL_NODE counter.
    - put_cpu_partial() function
        - See also: Slub Memory Allocator -7- (Object Unlocked) (http://jake.dothome.co.kr/slub-object-free) | Qc
- In lines 38~40 of code, if the per-cpu exceeds half of the S->cpu_partial value for the number of objects to manage, it exits the loop.
- Returns the free object of the first page moved from line 44 of code.

## acquire_slab()

mm/slub.c

```
 1  /*
 2   * Remove slab from the partial list, freeze it and
 3   * return the pointer to the freelist.
 4   *
 5   * Returns a list of objects or NULL if it fails.
 6   */
01  static inline void *acquire_slab(struct kmem_cache *s,
02                  struct kmem_cache_node *n, struct page *page,
03                  int mode, int *objects)
04  {
05          void *freelist;
06          unsigned long counters;
07          struct page new;
08
09          lockdep_assert_held(&n->list_lock);
10
11          /*
12           * Zap the freelist and set the frozen bit.
13           * The old freelist is the list of objects for the
14           * per cpu allocation list.
15           */
16          freelist = page->freelist;
17          counters = page->counters;
18          new.counters = counters;
19          *objects = new.objects - new.inuse;
20          if (mode) {
21                  new.inuse = page->objects;
22                  new.freelist = NULL;
23          } else {
24                  new.freelist = freelist;
25          }
26
27          VM_BUG_ON(new.frozen);
28          new.frozen = 1;
29
30          if (!__cmpxchg_double_slab(s, page,
31                          freelist, counters,
32                          new.freelist, new.counters,
33                          "acquire_slab"))
34                  return NULL;
35
```

```
36          remove_partial(n, page);
37          WARN_ON(!freelist);
38          return freelist;
39 }
```

In order to move one slap page from the partial list of per nodes to the CPU cache, it is detached by making it frozen. The value returned is the first free object, and if it fails, it returns null. According to the @mode, it works as follows:

- If @mode is true, the obtained slab page is obtained for the purpose of moving it to c->page. In order to make all the free objects on the slab page also in use, we substitute page->inuse with page->objects, and specify null for page->freelist.
- If the @mode is false, the obtained slab page is obtained for the purpose of moving it to a c->partial list.

## Slowpath Step 3 – Processing on Remote Nodes

### get_any_partial()

mm/slub.c

```
01 /*
02  * Get a page from somewhere. Search in increasing NUMA distances.
03  */
04 static void *get_any_partial(struct kmem_cache *s, gfp_t flags,
05              struct kmem_cache_cpu *c)
06 {
07 #ifdef CONFIG_NUMA
08          struct zonelist *zonelist;
09          struct zoneref *z;
10          struct zone *zone;
11          enum zone_type high_zoneidx = gfp_zone(flags);
12          void *object;
13          unsigned int cpuset_mems_cookie;
14
15          /*
16           * The defrag ratio allows a configuration of the tradeoffs betw
     een
17           * inter node defragmentation and node local allocations. A lowe
     r
18           * defrag_ratio increases the tendency to do local allocations
19           * instead of attempting to obtain partial slabs from other node
     s.
20           *
21           * If the defrag_ratio is set to 0 then kmalloc() always
22           * returns node local objects. If the ratio is higher then kmall
     oc()
23           * may return off node objects because partial slabs are obtaine
     d
24           * from other nodes and filled up.
25           *
26           * If /sys/kernel/slab/xx/remote_node_defrag_ratio is set to 100
27           * (which makes defrag_ratio = 1000) then every (well almost)
28           * allocation will first attempt to defrag slab caches on other
     nodes.
29           * This means scanning over all nodes to look for partial slabs
     which
30           * may be expensive if we do it every time we are trying to find
     a slab
31           * with available objects.
32           */
```

```
33          if (!s->remote_node_defrag_ratio ||
34                      get_cycles() % 1024 > s->remote_node_defrag_rati
    o)
35                  return NULL;
36
37          do {
38                  cpuset_mems_cookie = read_mems_allowed_begin();
39                  zonelist = node_zonelist(mempolicy_slab_node(), flags);
40                  for_each_zone_zonelist(zone, z, zonelist, high_zoneidx)
    {
41                          struct kmem_cache_node *n;
42
43                          n = get_node(s, zone_to_nid(zone));
44
45                          if (n && cpuset_zone_allowed(zone, flags) &&
46                                  n->nr_partial > s->min_partial)
    {
47                                  object = get_partial_node(s, n, c, flag
    s);
48                                  if (object) {
49                                          /*
50                                           * Don't check read_mems_allowed
    _retry()
51                                           * here - if mems_allowed was up
    dated in
52                                           * parallel, that was a harmless
    race
53                                           * between allocation and the cp
    uset
54                                           * update
55                                           */
56                                          return object;
57                                  }
58                          }
59                  }
60          } while (read_mems_allowed_retry(cpuset_mems_cookie));
61  #endif
62          return NULL;
63  }
```

Move the slab page from the partial list of remote nodes to the page and partial list of the per cpu cache in the cpu cache.

- On lines 39~42 of code, use the delay timer to perform the routine at the ratio of "/sys/kernel/slab/<slab >/remote_node_defrag_ratio=" (default=1000(99.8%), 1024=100%). If it is excluded from the probability, it returns null.
    - Delay Loop Timer (http://jake.dothome.co.kr/delay) | 문c
- In line 44~45 of the code, if there is a cpuset fluctuation, the cpuset sequence lock value is read **for retry**
- In line 46~47 of the code, traverse the zones of high_zoneidx or less from top to bottom in the zonelist of nodes found by mempolicy.
- In line 50~65 of the code, only if the node in the traversing zone is a node allowed by CPUset and the number of partial lists of nodes is greater than the minimum S->min_partial required by the cache, a certain amount of slab pages from the partial list of that node will be moved to the page and partial list of the per CPU cache. If successful, it returns the first free object.
- If there is a cpuset churn in line 67 of code, retry
- This is the case when the code is not processed to traverse the entire zonelist. Returns null as a failure.

# Atomic replacement of 2 values in slab

Architectures that provide an atomic API that allows you to change 86 (double) values at the same time on x64, arm390, s2, etc., will be handled with slightly faster performance, while those that don't will be replaced with interrupts blocked.

There are two related APIs:

- this_cpu_cmpxchg_double()
  - It is used to replace two values related to the per-cpu variable.
- cmpxchg_double()
  - It is used to replace two values related to a global variable.

Functions that atomically swap 2 values in a slab

- cmpxchg_double_slab()
  - If the architecture doesn't support an API that atomically swaps two values, disable interrupts in the function and handle them with bit_spin_lock.
- __cmpxchg_double_slab()
  - If the architecture doesn't support an API that atomically swaps two values, it will be handled with a bit_spin_lock in the function.
  - The interrupt must be disabled before the function is called.

## __cmpxchg_double_slab()

mm/slub.c

```
01  /* Interrupts must be disabled (for the fallback code to work right) */
02  static inline bool __cmpxchg_double_slab(struct kmem_cache *s, struct pa
    ge *page,
03                  void *freelist_old, unsigned long counters_old,
04                  void *freelist_new, unsigned long counters_new,
05                  const char *n)
06  {
07          VM_BUG_ON(!irqs_disabled());
08  #if defined(CONFIG_HAVE_CMPXCHG_DOUBLE) && \
09      defined(CONFIG_HAVE_ALIGNED_STRUCT_PAGE)
10          if (s->flags & __CMPXCHG_DOUBLE) {
11                  if (cmpxchg_double(&page->freelist, &page->counters,
12                                      freelist_old, counters_old,
13                                      freelist_new, counters_new))
14                          return 1;
15          } else
16  #endif
17          {
18                  slab_lock(page);
19                  if (page->freelist == freelist_old &&
20                                      page->counters == counters_old)
    {
21                          page->freelist = freelist_new;
```

```
22                        page->counters = counters_new;
23                        slab_unlock(page);
24                        return 1;
25                    }
26                slab_unlock(page);
27            }
28
29        cpu_relax();
30        stat(s, CMPXCHG_DOUBLE_FAIL);
31
32 #ifdef SLUB_DEBUG_CMPXCHG
33        pr_info("%s %s: cmpxchg double redo ", n, s->name);
34 #endif
35
36        return 0;
37 }
```

Compare page->freelist and counter to the old value and replace it with a new value if it is the same. In architectures that support CONFIG_HAVE_CMPXCHG_DOUBLE, the replacement of two word values is handled atomically. For unsupported architectures, use bit_spin_lock() to handle them.

- if @page->freelist == @freelist_old && @page->counters == @counters_old
    - @page->freelist    <— @freelist_new
    - @page->counters  <— @counters_new


## CONFIG_HAVE_CMPXCHG_DOUBLE

This is a kernel option used by architectures that support atomic APIs that replace 2 (double) values.

- x86, ARM64, and S390 architectures are supported.


The following figure shows two atomic operations: the fastpath used by the per CPU cache and the slowpath used by the per node.



(http://jake.dothome.co.kr/wp-content/uploads/2016/06/lock-for-slab-1b.png)

Here's how to atomically swap two values at once:

- Generic method of replacing two values after blocking an interrupt
- For fast performance, the architecture supports atomic API methods – such as the ARM64 architecture.

## this_cpu_cmpxchg_double()

include/linux/percpu-defs.h

```
1  #define this_cpu_cmpxchg_double(pcp1, pcp2, oval1, oval2, nval1, nval2) \
2          __pcpu_double_call_return_bool(this_cpu_cmpxchg_double_, pcp1, pcp2, oval1, oval2, nval1, nv
3  al2)
```

If the values of the two per-cpu variables pcp1 and pcp2 are the same as oval1 and oval2, change them to nval1 and nval2.

- if pcp1 == oval1 && pcp2 == oval2
  - pcp1 = nval1
  - pcp2 = nval2

## __pcpu_double_call_return_bool()

include/linux/percpu-defs.h

```
1  /*
2   * Special handling for cmpxchg_double.  cmpxchg_double is passed two
3   * percpu variables.  The first has to be aligned to a double word
4   * boundary and the second has to follow directly thereafter.
5   * We enforce this on all architectures even if they don't support
6   * a double cmpxchg instruction, since it's a cheap requirement, and it
7   * avoids breaking the requirement for architectures with the instructio
   n.
8   */
```

```
01  #define __pcpu_double_call_return_bool(stem, pcp1, pcp2, ...) \
02  ({ \
03          bool pdcrb_ret__; \
04          __verify_pcpu_ptr(&(pcp1)); \
05          BUILD_BUG_ON(sizeof(pcp1) != sizeof(pcp2)); \
06          VM_BUG_ON((unsigned long)(&(pcp1)) % (2 * sizeof(pcp1))); \
07          VM_BUG_ON((unsigned long)(&(pcp2)) != \
08                  (unsigned long)(&(pcp1)) + sizeof(pcp1)); \
09          switch(sizeof(pcp1)) { \
10          case 1: pdcrb_ret__ = stem##1(pcp1, pcp2, __VA_ARGS__); break; \
11          case 2: pdcrb_ret__ = stem##2(pcp1, pcp2, __VA_ARGS__); break; \
```

```
12          case 4: pdcrb_ret__ = stem##4(pcp1, pcp2, __VA_ARGS__); break;
   \
13          case 8: pdcrb_ret__ = stem##8(pcp1, pcp2, __VA_ARGS__); break;
   \
14          default:
   \
15                  __bad_size_call_parameter(); break;
   \
16          }
   \
17          pdcrb_ret__;
   \
18  })
```

For cmpxchg_double handling, call the function with a suffix of 1, 2, 4, or 8 to the function name that starts with @stem.

- e.g. @stem=this_cpu_cmpxchg_double_
    - this_cpu_cmpxchg_double_1(pcp1, pcp2, …)
    - this_cpu_cmpxchg_double_2(pcp1, pcp2, …)
    - this_cpu_cmpxchg_double_3(pcp1, pcp2, …)
    - this_cpu_cmpxchg_double_4(pcp1, pcp2, …)

### this_cpu_cmpxchg_double_8()

include/asm-generic/percpu.h

```
1  #ifndef this_cpu_cmpxchg_double_8
2  #define this_cpu_cmpxchg_double_8(pcp1, pcp2, oval1, oval2, nval1, nval2) \
3          this_cpu_generic_cmpxchg_double(pcp1, pcp2, oval1, oval2, nval1, nval2)
```

## Generic method using interrupt blocking

## this_cpu_generic_cmpxchg_double() – generic

include/asm-generic/percpu.h

```
01  #define this_cpu_generic_cmpxchg_double(pcp1, pcp2, oval1, oval2, nval1, nval2) \
02  ({
   \
03          int __ret;
   \
04          unsigned long __flags;
   \
05          raw_local_irq_save(__flags);
   \
06          __ret = raw_cpu_generic_cmpxchg_double(pcp1, pcp2,
   \
07                          oval1, oval2, nval1, nval2);
   \
08          raw_local_irq_restore(__flags);
   \
09          __ret;
   \
10  })
```

After masking the interrupt with a generic method, it compares the values of the two per-cpu variables and changes them to the new values.

**raw_cpu_generic_cmpxchg_double()**

include/asm-generic/percpu.h

```
01  #define raw_cpu_generic_cmpxchg_double(pcp1, pcp2, oval1, oval2, nval1,
    nval2) \
02  ({
    \
03          typeof(&(pcp1)) __p1 = raw_cpu_ptr(&(pcp1));
    \
04          typeof(&(pcp2)) __p2 = raw_cpu_ptr(&(pcp2));
    \
05          int __ret = 0;
    \
06          if (*__p1 == (oval1) && *__p2  == (oval2)) {
    \
07                  *__p1 = nval1;
    \
08                  *__p2 = nval2;
    \
09                  __ret = 1;
    \
10          }
    \
11          (__ret);
    \
12  })
```

Compare the values of the two per-cpu variables and change them to the new values.

# Atomic API methods supported by the architecture

## this_cpu_generic_cmpxchg_double() – ARM64

arch/arm64/include/asm/percpu.h

```
1  /*
2   * It would be nice to avoid the conditional call into the scheduler when
3   * re-enabling preemption for preemptible kernels, but doing that in a way
4   * which builds inside a module would mean messing directly with the preempt
5   * count. If you do this, peterz and tglx will hunt you down.
6   */
```

```
01  #define this_cpu_cmpxchg_double_8(ptr1, ptr2, o1, o2, n1, n2)
    \
02  ({
    \
03          int __ret;
    \
04          preempt_disable_notrace();
    \
05          __ret = cmpxchg_double_local(   raw_cpu_ptr(&(ptr1)),
    \
06                                          raw_cpu_ptr(&(ptr2)),
    \
```

```
07                                  o1, o2, n1, n2); \
08         preempt_enable_notrace(); \
09         __ret; \
10 })
```

After masking the preemption using the atomic instructions supported by the architecture, the values of the two per-cpu variables are compared and changed to the new values.

### cmpxchg_double_local() – ARM64

arch/arm64/include/asm/cmpxchg.h

```
1 #define cmpxchg_double_local(ptr1, ptr2, o1, o2, n1, n2) \
2 ({\
3         int __ret;\
4         __cmpxchg_double_check(ptr1, ptr2); \
5         __ret = !__cmpxchg_double((unsigned long)(o1), (unsigned long)(o
2), \
6                                   (unsigned long)(n1), (unsigned long)(n
2), \
7                                   ptr1); \
8         __ret; \
9 })
```

Use the atomic command, which is supported by the ARM64 architecture, to compare two unsigned long values and change them to the new values.

### __cmpxchg_double() Series – ARM64 – ll_sc Method

include/asm/atomic_ll_sc.h

```
01 #define __CMPXCHG_DBL(name, mb, rel, cl)                                \
02 __LL_SC_INLINE long                                                      \
03 __LL_SC_PREFIX(__cmpxchg_double##name(unsigned long old1,               \
04                                       unsigned long old2,               \
05                                       unsigned long new1,               \
06                                       unsigned long new2,               \
07                                       volatile void *ptr))              \
08 {                                                                        \
09         unsigned long tmp, ret;                                          \
10                                                                          \
11         asm volatile("// __cmpxchg_double" #name "\n"                   \
12         "       prfm    pstl1strm, %2\n"                                \
13         "1:     ldxp    %0, %1, %2\n"                                   \
14         "       eor     %0, %0, %3\n"                                   \
```

```
15          "           eor       %1, %1, %4\n"
   \
16          "           orr       %1, %0, %1\n"
   \
17          "           cbnz      %1, 2f\n"
   \
18          "           st" #rel "xp     %w0, %5, %6, %2\n"
   \
19          "           cbnz      %w0, 1b\n"
   \
20          "           " #mb "\n"
   \
21          "2:"
   \
22          : "=&r" (tmp), "=&r" (ret), "+Q" (*(unsigned long *)ptr)
   \
23          : "r" (old1), "r" (old2), "r" (new1), "r" (new2)
   \
24          : cl);
   \
25                                                                  \
26          return ret;
   \
27  }
   \
28  __LL_SC_EXPORT(__cmpxchg_double##name);
29
30  __CMPXCHG_DBL(    ,          ,  ,                )
31  __CMPXCHG_DBL(_mb, dmb ish, l, "memory")
```

The ARM64 architecture uses ll_sc atomic commands to compare two unsigned long values and change them to the new values.

## __cmpxchg_double() Series – ARM64 – LSE Method

```
01  #define __LL_SC_CMPXCHG_DBL(op) __LL_SC_CALL(__cmpxchg_double##op)
02
03  #define __CMPXCHG_DBL(name, mb, cl...)
    \
04  static inline long __cmpxchg_double##name(unsigned long old1,
    \
05                                            unsigned long old2,
    \
06                                            unsigned long new1,
    \
07                                            unsigned long new2,
    \
08                                            volatile void *ptr)
    \
09  {
    \
10          unsigned long oldval1 = old1;
    \
11          unsigned long oldval2 = old2;
    \
12          register unsigned long x0 asm ("x0") = old1;
    \
13          register unsigned long x1 asm ("x1") = old2;
    \
14          register unsigned long x2 asm ("x2") = new1;
    \
15          register unsigned long x3 asm ("x3") = new2;
    \
16          register unsigned long x4 asm ("x4") = (unsigned long)ptr;
    \
```

```
17                                                                          \
18          asm volatile(ARM64_LSE_ATOMIC_INSN(                            \
19          /* LL/SC */                                                     \
20          __LL_SC_CMPXCHG_DBL(name)                                       \
21          __nops(3),                                                      \
22          /* LSE atomics */                                               \
23          "       casp" #mb "\t%[old1], %[old2], %[new1], %[new2], %[v]\n"\
24          "       eor     %[old1], %[old1], %[oldval1]\n"                 \
25          "       eor     %[old2], %[old2], %[oldval2]\n"                 \
26          "       orr     %[old1], %[old1], %[old2]")                     \
27          : [old1] "+&r" (x0), [old2] "+&r" (x1),                         \
28            [v] "+Q" (*(unsigned long *)ptr)                             \
29          : [new1] "r" (x2), [new2] "r" (x3), [ptr] "r" (x4),            \
30            [oldval1] "r" (oldval1), [oldval2] "r" (oldval2)            \
31          : __LL_SC_CLOBBERS, ##cl);                                     \
32                                                                          \
33          return x0;                                                      \
34 }
35
36 __CMPXCHG_DBL(   ,    )
37 __CMPXCHG_DBL(_mb, al, "memory")
```

Using the lse-like atomic command, which is supported by the ARM64 architecture, two unsigned long values are compared and changed to the new values.

## cmpxchg_double() – generic

include/asm-generic/atomic-instrumented.h

```
1 #define cmpxchg_double(p1, p2, o1, o2, n1, n2)                         \
2 ({                                                                      \
3          typeof(p1) __ai_p1 = (p1);                                     \
4          kasan_check_write(__ai_p1, 2 * sizeof(*__ai_p1));              \
5          arch_cmpxchg_double(__ai_p1, (p2), (o1), (o2), (n1), (n2));    \
6 })
```

## cmpxchg_double() – ARM64

arch/arm64/include/asm/cmpxchg.h

```
1 #define cmpxchg_double(ptr1, ptr2, o1, o2, n1, n2) \
2 ({\
3          int __ret;\
```

```
 4            __cmpxchg_double_check(ptr1, ptr2); \
 5            __ret = !__cmpxchg_double_mb((unsigned long)(o1), (unsigned lon
  g)(o2), \
 6                                        (unsigned long)(n1), (unsigned lon
  g)(n2), \
 7                                        ptr1); \
 8            __ret; \
 9  })
```

# consultation

- Slab Memory Allocator -1- (Structure) (http://jake.dothome.co.kr/slub/) | Qc
- Slab Memory Allocator -2- (Initialize Cache) (http://jake.dothome.co.kr/kmem_cache_init) | Qc
- Slub Memory Allocator -3- (Create Cache) (http://jake.dothome.co.kr/slub-cache-create) | Qc
- Slub Memory Allocator -4- (Calculate Order) (http://jake.dothome.co.kr/slub-order) | Qc
- Slub Memory Allocator -5- | (http://jake.dothome.co.kr/slub-slub-alloc) Qc
- Slub Memory Allocator -6- (Assign Object) (http://jake.dothome.co.kr/slub-object-alloc) | Sentence C – Current post
- Slub Memory Allocator -7- (Object Unlocked) (http://jake.dothome.co.kr/slub-object-free) | Qc
- Slub Memory Allocator -8- (Drain/Flash Cache) (http://jake.dothome.co.kr/slub-drain-flush-cache) | Qc
- Slub Memory Allocator -9- (Cache Shrink) (http://jake.dothome.co.kr/slub-cache-shrink) | Qc
- Slub Memory Allocator -10- | (http://jake.dothome.co.kr/slub-slub-free) Qc
- Slub Memory Allocator -11- (Clear Cache (http://jake.dothome.co.kr/slub-cache-destroy)) | Qc
- Slub Memory Allocator -12- (Debugging Slub) (http://jake.dothome.co.kr/slub-debug) | Qc
- Slub Memory Allocator -13- (slabinfo) (http://jake.dothome.co.kr/slub-slabinfo) | 문c

# 12 thoughts to "Slub Memory Allocator -6- (Object 할당)"

**HYEYOO (HTTPS://HYEYOO.COM)**
2021-09-18 02:26 (http://jake.dothome.co.kr/slub-object-alloc/#comment-306028)

Hello, I have a question while studying, so I leave a comment. Slowpath – We have a function get_freelist from 1, and in this inline function we change the value of new, where new is a local variable.

new.inuse = page->objects;
new.frozen = freelist != NULL;

So why change the structure new? Is this accessible from outside the function get_freelist?
I think I've missed either the C language or the race condition.

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

2021-09-18 11:34 (http://jake.dothome.co.kr/slub-object-alloc/#comment-306030)

Hello?

The local variable new, which is created with a page structure, is being used temporarily, as you said.

However, new.counter is declared as a union, so if you check it, you'll understand it quickly.

new.counter = new.inuse | new.objects | new.frozen

```
참고) include/linux/mm_types.h

union {

void *s_mem; /* slab: first object */

unsigned long counters; /* SLUB */

struct { /* SLUB */

unsigned inuse:16;

unsigned objects:15;

unsigned frozen:1;

};

};
```

I appreciate it.

**HYEYOO (HTTPS://NAVER.COM-)**

2021-09-18 12:01 (http://jake.dothome.co.kr/slub-object-alloc/#comment-306032)

Thank you for your reply. But the gist of my question was what is the point of storing a value in new, a temporary variable used by get_freelist, and then not using it at all?

The function looks like this, but I didn't use it at all after changing new.

static inline void *get_freelist(struct kmem_cache *s, struct page *page)

{

struct page new;

unsigned long counters;

void *freelist;

```
do {
freelist = page->freelist;
counters = page->counters;

new.counters = counters;
VM_BUG_ON(!new.frozen);

new.inuse = page->objects;
new.frozen = freelist != NULL;

} while (!__cmpxchg_double_slab(s, page,
freelist, counters,
NULL, new.counters,
"get_freelist"));

return freelist;
}
```

RESPONSE (/SLUB-OBJECT-ALLOC/?REPLYTOCOM=306032#RESPOND)

---

**MOON YOUNG-IL (HTTP://JAKE.DOTHOME.CO.KR)**
2021-09-18 14:35 (http://jake.dothome.co.kr/slub-object-alloc/#comment-306033)

The cmpxchg_double_slab() command compares two values with two old values,
and then substitutes (replaces) two new values only if they are the same.
In the end, if
the values of page->freelist and page->counters are the same compared to freelist
and counters, then page->freelist = NULL, and page->counters = new.counters
values.

So we used the changed value of new. (new.inuse and new.frozen are reflected in
new.counters)
We didn't use all of the new structs, but we did use the new.counters section.

I appreciate it.

RESPONSE (/SLUB-OBJECT-ALLOC/?REPLYTOCOM=306033#RESPOND)

---

**HYEYOO (HTTPS://HYEYOO.COM)**
2021-09-18 16:47 (http://jake.dothome.co.kr/slub-object-alloc/#comment-306034)

Oh boy...!! I forgot that counters are a union of inuse/frozen/objects.
Not for bitwise, but for the use of unions! I appreciate it.

RESPONSE (/SLUB-OBJECT-ALLOC/?REPLYTOCOM=306034#RESPOND)

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

2021-09-18 18:58 (http://jake.dothome.co.kr/slub-object-alloc/#comment-306035)

Yes, fortunately, you found out right away. Have a nice day. ^^

RESPONSE (/SLUB-OBJECT-ALLOC/?REPLYTOCOM=306035#RESPOND)

**HYEYOO (HTTPS://HYEYOO.COM)**

2021-09-18 02:31 (http://jake.dothome.co.kr/slub-object-alloc/#comment-306029)

https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/commit/?
id=cc09ee80c3b1
(https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/commit/?
id=cc09ee80c3b1)
with the recent merge of PREEMPT_RT, slap and so on have changed a lot

RESPONSE (/SLUB-OBJECT-ALLOC/?REPLYTOCOM=306029#RESPOND)

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

2021-09-18 11:43 (http://jake.dothome.co.kr/slub-object-alloc/#comment-306031)

Thanks for the great information. I'll have to read it carefully.

Although PREEMPT_RT has been merged into the current kernel, some of the works
that have been done on the actual RT kernel have not yet been migrated for stability
and compatibility with existing kernels.
Thus, CONFIG_PREMPT_RT kernel options are slightly more preemtinable than
CONFIG_PREEMPT when enabled by individuals, but they are a little different from
the previous RT kernels.
This shortfall requires SoC vendors to do more work to release it.

I appreciate it.

RESPONSE (/SLUB-OBJECT-ALLOC/?REPLYTOCOM=306031#RESPOND)

**AS IT IS**

2021-11-06 19:40 (http://jake.dothome.co.kr/slub-object-alloc/#comment-306108)

Hello

I'm asking because I don't understand the get_freelist() function.

I compared freelist and page->freelist, counters and page->counters in do-while statements, and understood them to be syntax that exits the while statement if they are the same. One thing I'm curious about is when it happens that the two conditions are not the same. There is no syntax for re-assigning the values of the freelist and counters variables, and there doesn't seem to be any change to the page structure, so what are the cases where the two conditions are not the same?

RESPONSE (/SLUB-OBJECT-ALLOC/?REPLYTOCOM=306108#RESPOND)

**MOON YOUNG-IL (HTTP://JAKE.DOTHOME.CO.KR)**
2021-11-08 10:39 (http://jake.dothome.co.kr/slub-object-alloc/#comment-306111)

Hello?

In SMP systems, page->freelist can be contested when another cpu frees the page's slab object via the __slab_free() function.

I appreciate it.

RESPONSE (/SLUB-OBJECT-ALLOC/?REPLYTOCOM=306111#RESPOND)

**AS IT IS**
2021-11-08 21:32 (http://jake.dothome.co.kr/slub-object-alloc/#comment-306112)

Oh I see!

I've been looking at the SLUB allocator in my study lately, and I think the material you've put together really helps me a lot.
Thank you as always!

RESPONSE (/SLUB-OBJECT-ALLOC/?REPLYTOCOM=306112#RESPOND)

**MOON YOUNG-IL (HTTP://JAKE.DOTHOME.CO.KR)**
2021-11-09 09:01 (http://jake.dothome.co.kr/slub-object-alloc/#comment-306116)

I hope this helps you analyze. Have a nice day!

RESPONSE (/SLUB-OBJECT-ALLOC/?REPLYTOCOM=306116#RESPOND)

## LEAVE A COMMENT

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

Comments

name *

email *

Website

WRITE A COMMENT

Munc Blog (2015 ~ 2024)