


Linux 内核中创建 cache 节点由函数 `kmem_cache_create()` 实现。

该函数的执行流程：

- 1, 从全局 `cache_cache` 中获得 `cache` 结构, 因为全局 `cache_cache` 初始化对象的大小就是 `kmem_cache` 结构的大小, 所以返回的指针正好可以转换为 `cache` 结构; 调用 `kmem_cache_zalloc(&cache_cache, gfp)`;
- 2, 获得 slab 中碎片大小, 由函数 `calculate_slab_order()` 实现;
- 3, 计算并初始化 `cache` 的各种属性, 如果是外置式, 需要用 `kmem_find_general_cachep(slab_size, 0u)` 指定 `cachep->slabp_cache`, 用于存放 slab 对象和 `kmem_bufctl_t[]` 数组;
- 4, 设置每个 CPU 上得本地 `cache`, `setup_cpu_cache()`;
- 5, `cache` 创建完毕, 将其加入到全局 slab cache 链表中;

一、主实现

[cpp] [view plaincopyprint?](#)

```
1. /**
2.  * kmem_cache_create - Create a cache.
3.  * @name: A string which is used in /proc/slabinfo to identify this cache.
4.  * @size: The size of objects to be created in this cache.
5.  * @align: The required alignment for the objects.
6.  * @flags: SLAB flags
7.  * @ctor: A constructor for the objects.
8.  *
9.  * Returns a ptr to the cache on success, NULL on failure.
10. * Cannot be called within a int, but can be interrup.
11. * The @ctor is run when new pages are allocated by the cache.
12. *
13. * @name must be valid until the cache is destroyed. This implies that
14. * the module calling this has to destroy the cache before getting unloaded.
15. * Note that kmem_cache_name() is not guaranteed to return the same pointer,
16. * therefore applications must manage it themselves.
17. *
18. * The flags are
19. *
20. * %SLAB_POISON - Poison the slab with a known test pattern (a5a5a5a5)
21. * to catch references to uninitialised memory.
22. *
23. * %SLAB_RED_ZONE - Insert 'Red' zones around the allocated memory to check
```

```

24. * for buffer overruns.
25. *
26. * %SLAB_HWCACHE_ALIGN - Align the objects in this cache to a hardware
27. * cacheline. This can be beneficial if you're counting cycles as closely
28. * as davem.
29. */
30. /*创建 slab 系统顶层的 cache 节点。创建完成后，cache
31. 里并没有任何 slab 以及对象，只有当分配对象
32. ，并且 cache 中没有空闲对象时，才会创建新的 slab。*/
33. struct kmem_cache *
34. kmem_cache_create (const char *name, size_t size, size_t align,
35. unsigned long flags, void (*ctor)(void *))
36. {
37.     size_t left_over, slab_size, ralign;
38.     struct kmem_cache *cachep = NULL, *pc;
39.     gfp_t gfp;
40.
41.     /*
42.      * Sanity checks... these are all serious usage bugs.
43.      **/
44.     if (!name || in_interrupt() || (size < BYTES_PER_WORD) ||
45.         size > KMALLOC_MAX_SIZE) {
46.         printk(KERN_ERR "%s: Early error in slab %s\n", __func__,
47.             name);
48.         BUG();
49.     }
50.
51.     /*
52.      * We use cache_chain_mutex to ensure a consistent view of
53.      * cpu_online_mask as well. Please see cpuup_callback
54.      */
55.     /* slab 分配器是否已经初始化好，如果是内核启动阶段
56.      ，则只有一个 cpu 执行 slab 分配器的初始化动作，无需加锁，否则需要加锁 */
57.     if (slab_is_available()) {
58.         get_online_cpus();
59.         mutex_lock(&cache_chain_mutex);
60.     }
61.     /* 遍历 cache 链，做些校验工作 */
62.     list_for_each_entry(pc, &cache_chain, next) {
63.         char tmp;
64.         int res;
65.
66.         /*

```

```

67.      * This happens when the module gets unloaded and doesn't
68.      * destroy its slab cache and no-one else reuses the vmalloc
69.      * area of the module. Print a warning.
70.      */
71.      /* 检查 cache 链表中的 cache 是否都有名字 */
72.      res = probe_kernel_address(pc->name, tmp);
73.      if (res) { /* 没有名字, 报错 */
74.          printk(KERN_ERR
75.                  "SLAB: cache with size %d has lost its name\n",
76.                  pc->buffer_size);
77.          continue;
78.      }
79.      /* 检查 cache 链表中是否已经存在相同名字的 cache */
80.      if (!strcmp(pc->name, name)) {
81.          printk(KERN_ERR
82.                  "kmem_cache_create: duplicate cache %s\n", name);
83.          dump_stack();
84.          goto oops;
85.      }
86.  }
87.
88. #if DEBUG
89.  WARN_ON(strchr(name, ' ')); /* It confuses parsers */
90. #if FORCED_DEBUG
91.  /*
92.   * Enable redzoning and last user accounting, except for caches with
93.   * large objects, if the increased size would increase the object size
94.   * above the next power of two: caches with object sizes just above a
95.   * power of two have a significant amount of internal fragmentation.
96.   */
97.  if (size < 4096 || fls(size - 1) == fls(size - 1 + REDZONE_ALIGN +
98.          2 * sizeof(unsigned long long)))
99.      flags |= SLAB_RED_ZONE | SLAB_STORE_USER;
100.  if (!(flags & SLAB_DESTROY_BY_RCU))
101.      flags |= SLAB_POISON;
102. #endif
103.  if (flags & SLAB_DESTROY_BY_RCU)
104.      BUG_ON(flags & SLAB_POISON);
105. #endif
106.  /*
107.   * Always checks flags, a caller might be expecting debug support which
108.   * isn't available.
109.   */

```

```

110. BUG_ON(flags & ~CREATE_MASK);
111.
112. /*
113.  * Check that size is in terms of words. This is needed to avoid
114.  * unaligned accesses for some archs when redzoning is used, and makes
115.  * sure any on-slab bufctl's are also correctly aligned.
116.  */
117. if (size & (BYTES_PER_WORD - 1)) {
118.     size += (BYTES_PER_WORD - 1);
119.     size &= ~(BYTES_PER_WORD - 1);
120. }
121.
122. /* calculate the final buffer alignment: */
123.
124. /* 1) arch recommendation: can be overridden for debug */
125. if (flags & SLAB_HWCACHE_ALIGN) {
126.     /*
127.      * Default alignment: as specified by the arch code. Except if
128.      * an object is really small, then squeeze multiple objects into
129.      * one cacheline.
130.      */
131.     ralign = cache_line_size();
132.     while (size <= ralign / 2)
133.         ralign /= 2;
134. } else {
135.     ralign = BYTES_PER_WORD;
136. }
137.
138. /*
139.  * Redzoning and user store require word alignment or possibly larger.
140.  * Note this will be overridden by architecture or caller mandated
141.  * alignment if either is greater than BYTES_PER_WORD.
142.  */
143. if (flags & SLAB_STORE_USER)
144.     ralign = BYTES_PER_WORD;
145.
146. if (flags & SLAB_RED_ZONE) {
147.     ralign = REDZONE_ALIGN;
148.     /* If redzoning, ensure that the second redzone is suitably
149.      * aligned, by adjusting the object size accordingly. */
150.     size += REDZONE_ALIGN - 1;
151.     size &= ~(REDZONE_ALIGN - 1);
152. }

```

```

153.
154. /* 2) arch mandated alignment */
155. if (ralign < ARCH_SLAB_MINALIGN) {
156.     ralign = ARCH_SLAB_MINALIGN;
157. }
158. /* 3) caller mandated alignment */
159. if (ralign < align) {
160.     ralign = align;
161. }
162. /* disable debug if necessary */
163. if (ralign > __alignof__(unsigned long long))
164.     flags &= ~(SLAB_RED_ZONE | SLAB_STORE_USER);
165. /*
166.  * 4) Store it.
167.  */
168. align = ralign;
169. /* slab 分配器是否已经可用 */
170. if (slab_is_available())
171.     gfp = GFP_KERNEL;
172. else
173.     /* slab 初始化好之前，不允许阻塞，且只能在低端内存区分配 */
174.     gfp = GFP_NOWAIT;
175.
176. /* Get cache's description obj. */
177. /* 获得 struct kmem_cache 对象，为什么能从 cache 中获得的对象是
178. kmem_cache 结构呢，因为这里的全局变量 cache_cache 的对象大小
179. 就是 kmem_cache 结构大小*/
180. cachep = kmem_cache_zalloc(&cache_cache, gfp);
181. if (!cachep)
182.     goto oops;
183.
184. #if DEBUG
185.     cachep->obj_size = size;
186.
187.     /*
188.      * Both debugging options require word-alignment which is calculated
189.      * into align above.
190.      */
191.     if (flags & SLAB_RED_ZONE) {
192.         /* add space for red zone words */
193.         cachep->obj_offset += sizeof(unsigned long long);
194.         size += 2 * sizeof(unsigned long long);
195.     }

```

```

196. if (flags & SLAB_STORE_USER) {
197.     /* user store requires one word storage behind the end of
198.      * the real object. But if the second red zone needs to be
199.      * aligned to 64 bits, we must allow that much space.
200.      */
201.     if (flags & SLAB_RED_ZONE)
202.         size += REDZONE_ALIGN;
203.     else
204.         size += BYTES_PER_WORD;
205. }
206. #if FORCED_DEBUG && defined(CONFIG_DEBUG_PAGEALLOC)
207.     if (size >= malloc_sizes[INDEX_L3 + 1].cs_size
208.         && cachep->obj_size > cache_line_size() && size < PAGE_SIZE) {
209.         cachep->obj_offset += PAGE_SIZE - size;
210.         size = PAGE_SIZE;
211.     }
212. #endif
213. #endif
214.
215. /*
216.  * Determine if the slab management is 'on' or 'off' slab.
217.  * (bootstrapping cannot cope with offslab caches so don't do
218.  * it too early on.)
219.  */
220. /* 确定 slab 管理对象的存储方式：内置还是外置
221.  * 。通常，当对象大于等于 512 时，使用外置方式
222.  * 。初始化阶段采用内置式。
223.  * slab_early_init 参见 kmem_cache_init 函数 */
224. if ((size >= (PAGE_SIZE >> 3)) && !slab_early_init)
225.     /*
226.      * Size is large, assume best to place the slab management obj
227.      * off-slab (should allow better packing of objs).
228.      */
229.     flags |= CFLGS_OFF_SLAB;
230.
231. size = ALIGN(size, align);
232. /* 获得 slab 中碎片的大小 */
233. left_over = calculate_slab_order(cachep, size, align, flags);
234. /* cachep->num 为该 cache 中每个 slab 的对象数，为 0，表示为该对象创建 cache 失败 */
235. if (!cachep->num) {
236.     printk(KERN_ERR
237.         "kmem_cache_create: couldn't create cache %s.\n", name);
238.     kmem_cache_free(&cache_cache, cachep);

```

```

239.     cachep = NULL;
240.     goto oops;
241. }
242. /* 计算 slab 管理对象的大小，包括 struct slab 对象和 kmem_bufctl_t 数组 */
243. slab_size = ALIGN(cachep->num * sizeof(kmem_bufctl_t)
244.     + sizeof(struct slab), align);
245.
246. /*
247.  * If the slab has been placed off-slab, and we have enough space then
248.  * move it on-slab. This is at the expense of any extra colouring.
249.  */
250.
251. /* 如果这是一个外置式 slab，并且碎片大小大于 slab 管理对象的大小
252.  * ，则可将 slab 管理对象移到 slab 中，改造成一个内置式 slab */
253. if (flags & CFLGS_OFF_SLAB && left_over >= slab_size) {
254.     /* 除去 off-slab 标志位 */
255.     flags &= ~CFLGS_OFF_SLAB;
256.     /* 更新碎片大小 */
257.     left_over -= slab_size;
258. }
259.
260. if (flags & CFLGS_OFF_SLAB) {
261.     /* really off slab. No need for manual alignment */
262.     /* align 是针对 slab 对象的，如果 slab 管理对象是外置存储
263.     * ，自然不会像内置那样影响到后面 slab 对象的存储位置
264.     * ，也就不需要对齐了 */
265.     slab_size =
266.         cachep->num * sizeof(kmem_bufctl_t) + sizeof(struct slab);
267.
268. #ifdef CONFIG_PAGE_POISONING
269.     /* If we're going to use the generic kernel_map_pages()
270.     * poisoning, then it's going to smash the contents of
271.     * the redzone and userword anyhow, so switch them off.
272.     */
273.     if (size % PAGE_SIZE == 0 && flags & SLAB_POISON)
274.         flags &= ~(SLAB_RED_ZONE | SLAB_STORE_USER);
275. #endif
276. }
277. /* cache 的着色块的单位大小 */
278. cachep->colour_off = cache_line_size();
279. /* Offset must be a multiple of the alignment. */
280. /* 着色块大小必须是对象要求对齐方式的倍数 */
281. if (cachep->colour_off < align)

```

```

282.     cachep->colour_off = align;
283.     /* 计算碎片区需要多少个着色快 */
284.     cachep->colour = left_over / cachep->colour_off;
285.     /* slab 管理对象的大小 */
286.     cachep->slab_size = slab_size;
287.     cachep->flags = flags;
288.     cachep->gfpflags = 0;
289.     if (CONFIG_ZONE_DMA_FLAG && (flags & SLAB_CACHE_DMA))
290.         cachep->gfpflags |= GFP_DMA;
291.     /* slab 对象的大小 */
292.     cachep->buffer_size = size;
293.     /* 计算对象在 slab 中索引时用, 参见 obj_to_index 函数 */
294.     cachep->reciprocal_buffer_size = reciprocal_value(size);
295.
296.     if (flags & CFLGS_OFF_SLAB) {
297.         /* 分配一个 slab 管理区域对象, 保存在 slabp_cache 中,
298.            这个函数传入的大小为 slab_size, 也就是分配 slab_size 大小的 cache
299.            ,在 slab 创建的时候如果是外置式, 那么需要从分配的这里面
300.            分配出 slab 对象, 剩下的空间放 kmem_bufctl_t[] 数组,
301.            如果是内置式的 slab, 此指针为空 */
302.         cachep->slabp_cache = kmem_find_general_cachep(slab_size, 0u);
303.         /*
304.          * This is a possibility for one of the malloc_sizes caches.
305.          * But since we go off slab only for object size greater than
306.          * PAGE_SIZE/8, and malloc_sizes gets created in ascending order,
307.          * this should not happen at all.
308.          * But leave a BUG_ON for some lucky dude.
309.          */
310.         BUG_ON(ZERO_OR_NULL_PTR(cachep->slabp_cache));
311.     }
312.     cachep->ctor = ctor;
313.     cachep->name = name;
314.     /* 设置每个 cpu 上的 local cache */
315.     if (setup_cpu_cache(cachep, gfp)) {
316.         __kmem_cache_destroy(cachep);
317.         cachep = NULL;
318.         goto oops;
319.     }
320.
321.     /* cache setup completed, link it into the list */
322.     /* cache 创建完毕, 将其加入到全局 slab cache 链表中 */
323.     list_add(&cachep->next, &cache_chain);
324.oops:

```



```

325. if (!cachep && (flags & SLAB_PANIC))
326.     panic("kmem_cache_create(): failed to create slab `%s`\n",
327.           name);
328. if (slab_is_available()) {
329.     mutex_unlock(&cache_chain_mutex);
330.     put_online_cpus();
331. }
332. return cachep;
333.}

```

其中，cache_cache

[cpp] [view plaincopyprint?](#)

```

1. /* internal cache of cache description objs */
2. static struct kmem_cache cache_cache = {
3.     .batchcount = 1,
4.     .limit = BOOT_CPUCACHE_ENTRIES,
5.     .shared = 1,
6.     .buffer_size = sizeof(struct kmem_cache),/*大小为 cache 结构，难怪名称为 cache_cache*/
7.     .name = "kmem_cache",
8. };

```

二、计算 slab 碎片大小

[cpp] [view plaincopyprint?](#)

```

1. /**
2.  * calculate_slab_order - calculate size (page order) of slabs
3.  * @cachep: pointer to the cache that is being created
4.  * @size: size of objects to be created in this cache.
5.  * @align: required alignment for the objects.
6.  * @flags: slab allocation flags
7.  *
8.  * Also calculates the number of objects per slab.
9.  *
10. * This could be made much more intelligent. For now, try to avoid using
11. * high order pages for slabs. When the gfp() functions are more friendly
12. * towards high-order requests, this should be changed.
13. */
14. /*计算 slab 由几个页面组成，同时计算每个 slab 中有多少对象*/
15. static size_t calculate_slab_order(struct kmem_cache *cachep,
16.     size_t size, size_t align, unsigned long flags)
17. {
18.     unsigned long offslab_limit;
19.     size_t left_over = 0;

```

```

20. int gfporder;
21.
22. for (gfporder = 0; gfporder <= KMALLOC_MAX_ORDER; gfporder++) {
23.     unsigned int num;
24.     size_t remainder;
25.     /* 计算 slab 中对象数 */
26.     cache_estimate(gfporder, size, align, flags, &remainder, &num);
27.     /* 对象数为 0，表示此 order 下，一个对象都放不下，检查下一 order */
28.     if (!num)
29.         continue;
30.
31.     if (flags & CFLGS_OFF_SLAB) {
32.         /*
33.          * Max number of objs-per-slab for caches which
34.          * use off-slab slabs. Needed to avoid a possible
35.          * looping condition in cache_grow().
36.          */
37.         /* 创建一个外置式 slab 时，要相应分配该 slab 的管理对象
38.          * ，包含 struct slab 对象和 kmem_bufctl_t 数组，分配管理对象的流程就是分配普通对象
39.          * 的流程
40.          * ，再来看一下分配对象的流程：
41.          * kmem_cache_alloc->__cache_alloc->__do_cache_alloc->__cache_alloc->cache
42.          * _alloc_refill->cache_grow->alloc_slabmgmt->kmem_cache_alloc_node->kmem_cache_a
43.          * lloc
44.          * 可以看出这里可能存在一个循环，循环的关键在于 alloc_slabmgmt 函数
45.          * ，当 slab 管理对象是 off-slab 方式时，就形成了循环
46.          * 。那么什么时候 slab 管理对象会采用外置式 slab 呢？显然当其管理的 slab 中对象很多
47.          * ，从而 kmem_bufctl_t 数组很大，致使整个管理对象也很大，此时才会形成循环
48.          * 。故需要对 kmem_bufctl_t 的数目做限制，下面的算法是很粗略的，既然对象大小为
49.          * size 时
50.          * ，是外置式 slab，那么我们假设管理对象的大小也是 size，计算出 kmem_bufctl_t 数组
51.          * 的大小
52.          * ，即此大小的 kmem_bufctl_t 数组一定会造成管理对象是外置式 slab。之所以说粗略
53.          * ，是指数组大小小于这个限制时，也不能确保管理对象一定是内置式 slab。但这也不会引
54.          * 发错误
55.          * ，因为还有一个 slab_break_gfp_order 阀门来控制每个 slab 所占页面数，通常其值为
56.          * 1，即每个 slab 最多两个页面
57.          * ，外置式 slab 存放的都是大于 512 的大对象，所以
58.          * slab 中不会有太多的大对象，kmem_bufctl_t 数组也不会很大，粗略判断一下就足够了。
59.          */
60.         offslab_limit = size - sizeof(struct slab);
61.         offslab_limit /= sizeof(kmem_bufctl_t);
62.         /* 对象数目大于限制，跳出循环，不再尝试更大的 order

```

```

56.     , 避免 slab 中对象数目过多
57.     , 此时计算的对象数也是有效的, 循环一次没什么 */
58.     if (num > offslab_limit)
59.         break;
60. }
61.
62. /* Found something acceptable - save it away */
63. /* 每个 slab 中的对象数 */
64. cachep->num = num;
65. /* slab 的 order , 即由几个页面组成 */
66. cachep->gfporder = gfporder;
67. /* slab 中剩余空间 ( 碎片 ) 的大小 */
68. left_over = remainder;
69.
70. /*
71.  * A VFS-reclaimable slab tends to have most allocations
72.  * as GFP_NOFS and we really don't want to have to be allocating
73.  * higher-order pages when we are unable to shrink dcache.
74.  */
75. /* SLAB_RECLAIM_ACCOUNT 表示此 slab 所占页面为可回收的
76.  , 当内核检测是否有足够的页面满足用户态的需求时
77.  , 此类页面将被计算在内, 通过调用
78.  kmem_freepages()函数可以释放分配给 slab 的页框。由于是可回收的
79.  , 所以不需要做后面的碎片检测了 */
80. if (flags & SLAB_RECLAIM_ACCOUNT)
81.     break;
82.
83. /*
84.  * Large number of objects is good, but very large slabs are
85.  * currently bad for the gfp()s.
86.  */
87. /* slab_break_gfp_order 为 slab 所占页面的阀门, 超过这个阀门时
88.  , 无论碎片大小, 都不再检测更高的 order 了 */
89. if (gfporder >= slab_break_gfp_order)
90.     break;
91.
92. /*
93.  * Acceptable internal fragmentation?
94.  */
95. /* slab 所占页面的大小是碎片大小的 8 倍以上
96.  , 页面利用率较高, 可以接受这样的 order */
97. if (left_over * 8 <= (PAGE_SIZE << gfporder))
98.     break;

```

```
99. }
100. /* 返回碎片大小 */
101. return left_over;
102.}
```

三、查找指定大小 cache

[cpp] [view plaincopyprint?](#)

```
1. /*在 general cache 中分配一个 struct kmem_cache 对象。直接调用__find_general_cache。*/
2. static struct kmem_cache *kmem_find_general_cache(size_t size, gfp_t gfpflags)
3. {
4.     return __find_general_cache(size, gfpflags);
5. }
```

[cpp] [view plaincopyprint?](#)

```
1. static inline struct kmem_cache *__find_general_cache(size_t size,
2.                                                         gfp_t gfpflags)
3. {
4.     struct cache_sizes *csizes = malloc_sizes;
5.
6.     #if DEBUG
7.         /* This happens if someone tries to call
8.          * kmem_cache_create(), or __kmalloc(), before
9.          * the generic caches are initialized.
10.         */
11.         BUG_ON(malloc_sizes[INDEX_AC].cs_cache == NULL);
12. #endif
13.     if (!size)
14.         return ZERO_SIZE_PTR;
15.     /* 找到合适的 malloc size */
16.     while (size > csizes->cs_size)
17.         csizes++;
18.
19.     /*
20.      * Really subtle: The last entry with cs->cs_size==ULONG_MAX
21.      * has cs_{dma,}cache==NULL. Thus no special case
22.      * for large kmalloc calls required.
23.      */
24. #ifdef CONFIG_ZONE_DMA
25.     if (unlikely(gfpflags & GFP_DMA))
26.         return csizes->cs_dmacache;
27. #endif
28.     /* 返回该大小级别的 cache */
```

```
29. return csizep->cs_cachep;
30.}
```

四、设置 CPU 本地 cache

[cpp] [view plaincopyprint?](#)

```
1. /*配置 local cache 和 slab 三链。*/
2. static int __init_refok setup_cpu_cache(struct kmem_cache *cachep, gfp_t gfp)
3. {
4.     /* general cache 初始化完毕，配置每个 cpu 的 local cache */
5.     if (g_cpucache_up == FULL)
6.         return enable_cpucache(cachep, gfp);
7.     /* 此时处于系统初始化阶段，g_cpucache_up 记录 general cache 初始化的进度
8.      *，比如 PARTIAL_AC 表示 struct array_cache 所在的 cache 已经创建，
9.      * PARTIAL_L3 表示 struct kmem_list3 所在的 cache 已经创建
10.    *，注意创建这两个 cache 的先后顺序
11.    *。在初始化阶段只需配置主 cpu 的 local cache 和 slab 三链 */
12.     if (g_cpucache_up == NONE) {
13.         /*
14.          * Note: the first kmem_cache_create must create the cache
15.          * that's used by kmalloc(24), otherwise the creation of
16.          * further caches will BUG().
17.          */
18.         /* 初始化阶段创建 struct array_cache 所在 cache 时进入这个流程
19.          *，此时 struct array_cache 所在的 general cache 还未创建
20.          *，只能使用静态分配的全局变量 initarray_generic 表示的 local cache */
21.         cachep->array[smp_processor_id()] = &initarray_generic.cache;
22.
23.         /*
24.          * If the cache that's used by kmalloc(sizeof(kmem_list3)) is
25.          * the first cache, then we need to set up all its list3s,
26.          * otherwise the creation of further caches will BUG().
27.          */
28.         /* 创建 struct kmem_list3 所在的 cache 是在 struct array_cache 所在 cache 之后
29.          *，所以此时 struct kmem_list3 所在的
30.          * cache 也一定没有创建，也需要使用全局变量 */
31.         set_up_list3s(cachep, SIZE_AC);
32.         /* 执行到这 struct array_cache 所在的 cache 创建完毕
33.          *，如果 struct kmem_list3 和 struct array_cache 位于同一个 general cache 中
34.          *，不会再重复创建了
35.          *，g_cpucache_up 表示的进度更进一步 */
36.         if (INDEX_AC == INDEX_L3)
37.             g_cpucache_up = PARTIAL_L3;
38.         else
```

```

39.     g_cpucache_up = PARTIAL_AC;
40. } else {
41.     /* g_cpucache_up 至少为 PARTIAL_AC 时进入这个流程，struct array_cache 所在的
42.     general cache 已经建立起来，可以通过 kmalloc 分配了 */
43.     cachep->array[smp_processor_id()] =
44.         kmalloc(sizeof(struct arraycache_init), gfp);
45.
46.     if (g_cpucache_up == PARTIAL_AC) {
47.         /* struct kmem_list3 所在 cache 仍未创建完毕，还需使用全局的 slab 三链 */
48.         set_up_list3s(cachep, SIZE_L3);
49.         /* 后面将会分析 kmem_cache_init 函数，只有创建 struct kmem_list3 所在
50.         cache 时才会进入此流程，上面的代码执行完，struct kmem_list3 所在
51.         cache 也就创建完毕可以使用了，更新 g_cpucache_up */
52.         g_cpucache_up = PARTIAL_L3;
53.     } else {
54.         int node;
55.         for_each_online_node(node) {
56.             cachep->nodelists[node] = /* 通过 kmalloc 分配 struct kmem_list3 对象 */
57.                 kmalloc_node(sizeof(struct kmem_list3),
58.                             gfp, node);
59.             BUG_ON(!cachep->nodelists[node]);
60.             /* 初始化 slab 三链 */
61.             kmem_list3_init(cachep->nodelists[node]);
62.         }
63.     }
64. }
65. /* 设置回收时间 */
66. cachep->nodelists[numa_node_id()]->next_reap =
67.     jiffies + REAPTIMEOUT_LIST3 +
68.     ((unsigned long)cachep) % REAPTIMEOUT_LIST3;
69.
70. cpu_cache_get(cachep)->avail = 0;
71. cpu_cache_get(cachep)->limit = BOOT_CPUCACHE_ENTRIES;
72. cpu_cache_get(cachep)->batchcount = 1;
73. cpu_cache_get(cachep)->touched = 0;
74. cachep->batchcount = 1;
75. cachep->limit = BOOT_CPUCACHE_ENTRIES;
76. return 0;
77.}

```