# 代码导读5 - slab分配机制

## 笨叔叔

# 课程目标

- 为什么需要有slab分配机制？
- 使用slab机制有什么好处？或者说不使用slab有什么坏处？
- 怎么理解slab机制里面的cache，slab，object概念？
- Slab机制的着色区怎么理解？
- 一个slab可以装载多少个对象？需要多少个物理页面？怎么计算的？
- 使用slab机制分配一个字节的对象时候，实际分配是1个字节吗？
- 使用**kmem_cache_alloc()**函数获取空闲对象时候，**slab**什么时候为这些空闲的对象们分配物理页面？
- 我使用**kmem_cache_create()**函数创建了一个自己的slab cache，为啥我在slabinfo中没找到呢？
- slab调试的基本技巧

# SLAB机制API

#创建slab描述符

struct kmem_cache *

**kmem_cache_create**(const char *name, size_t size, size_t align,

unsigned long flags, void (*ctor)(void *))

#释放slab描述符

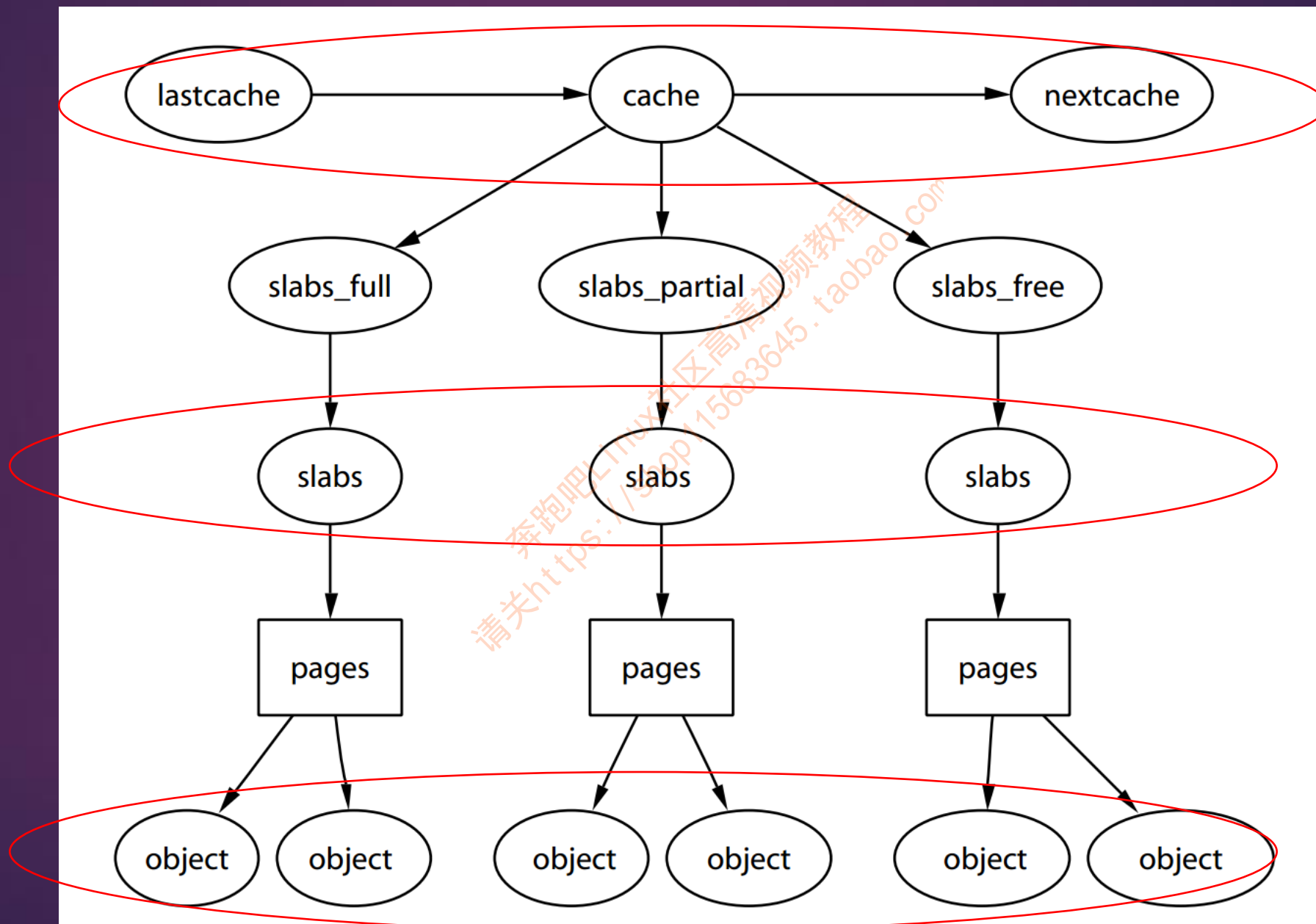void kmem_cache_destroy(struct kmem_cache *s)

#分配缓存对象

void **kmem_cache_alloc**(struct kmem_cache *, gfp_t flags);

#释放缓存对象

void kmem_cache_free(struct kmem_cache *, void *);

开源精神  奔跑不息

# 酒吧 如何存放啤酒的



存酒的小仓库

存酒的箱子

啤酒

# 酒吧小仓库是如何描述的呢？

```c
 6 /*
 7  * Definitions unique to the original Linux SLAB allocator.
 8  */
 9
10 struct kmem_cache {
11         struct array_cache __percpu *cpu_cache;
12
13 /* 1) Cache tunables. Protected by slab_mutex */
14         unsigned int batchcount;
15         unsigned int limit;
16         unsigned int shared;
17
18         unsigned int size;
19         struct reciprocal_value reciprocal_buffer_size;
20 /* 2) touched by every alloc & free from the backend */
21
22         unsigned int flags;             /* constant flags */
23         unsigned int num;               /* # of objs per slab */
24
25 /* 3) cache_grow/shrink */
26         /* order of pgs per slab (2^n) */
27         unsigned int gfporder;
28
29         /* force GFP flags, e.g. GFP_DMA */
30         gfp_t allocflags;
31
32         size_t colour;                  /* cache colouring range */
33         unsigned int colour_off;        /* colour offset */
34         struct kmem_cache *freelist_cache;
35         unsigned int freelist_size;
```

奔跑吧

# 重点分析函数1 – kmem_cache_create

```
361 struct kmem_cache *
362 kmem_cache_create(const char *name, size_t size, size_t align,
363                   unsigned long flags, void (*ctor)(void *))
364 {
365         struct kmem_cache *s;
366         const char *cache_name;
367         int err;
368
```
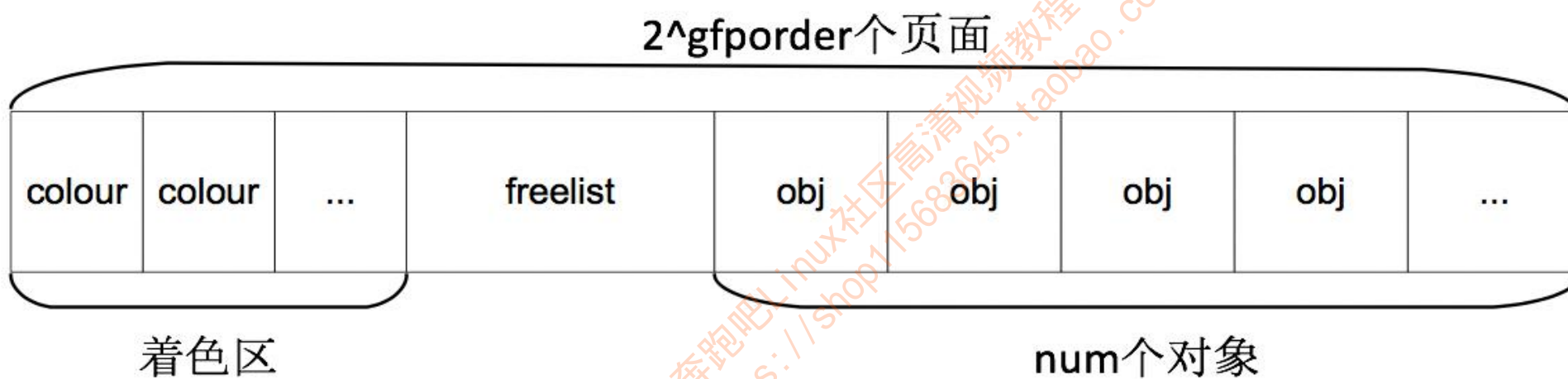
Name: 仓库的名称
Size：对象的大小
Align：对象对齐的要求
Flags：cache分配需求集合
Ctor：对象额外的构造函数

# Slab构造结构图

# 一个slab可以存放几个对象？

```
1923 static size_t calculate_slab_order(struct kmem_cache *cachep,
1924                         size_t size, size_t align, unsigned long flags)
1925 {
1926         unsigned long offslab_limit;
1927         size_t left_over = 0;
1928         int gfporder;
1929
1930         for (gfporder = 0; gfporder <= KMALLOC_MAX_ORDER; gfporder++) {
1931                 unsigned int num;
1932                 size_t remainder;
1933
1934                 cache_estimate(gfporder, size, align, flags, &remainder, &num);
1935                 if (!num)
1936                         continue;
1937
1938                 /* Can't handle number of objects more than SLAB_OBJ_MAX_NUM */
1939                 if (num > SLAB_OBJ_MAX_NUM)
1940                         break;
1941
```

计算公式：

$$obj\_num = slab\text{-}szie \: / \: (obj\_size + sizeof(freelist\_idx\_t))$$

奔跑吧Linux社区出品      开源精神 奔跑不息

# kmem_cache_create创建完后小仓库长啥样？

在ARM Vexpress平台上创建名为"figo_object"的slab描述符，大小为20Byte，align为8Byte，flags为0，假设L1 cache line大小为16Byte，其slab描述符相关成员的计算结果如下：

```
struct kmem_cache *cachep {
.array_cache = {
        .avail =0,
        .limit = 120,
        .batchmount = 60,
        .touched = 0,
 },
.batchount = 60,
.limit = 120,
.shared = 8,
.size = 24,
.flags = 0,
.num = 163,
.gfporder = 0,
.colour = 1,
.colour_off = 16,
.freelist_size = 168,
.name = "figo_object",
.object_size = 20,
.align =8,
.kmem_cache_node = {
    .free_object = 0,
    .free_limit = 283,
    .shared = {
        .avail =0,
        .limit = 480,
    },
 },
}
```

奔跑

# kmem_cache_alloc函数分析

```c
3385 /**
3386  * kmem_cache_alloc - Allocate an object
3387  * @cachep: The cache to allocate from.
3388  * @flags: See kmalloc().
3389  *
3390  * Allocate an object from this cache.  The flags are only relevant
3391  * if the cache has no available objects.
3392  */
3393 void *kmem_cache_alloc(struct kmem_cache *cachep, gfp_t flags)
3394 {
3395         void *ret = slab_alloc(cachep, flags, _RET_IP_);
3396
3397         trace_kmem_cache_alloc(_RET_IP_, ret,
3398                                cachep->object_size, cachep->size, flags);
3399
3400         return ret;
3401 }
3402 EXPORT_SYMBOL(kmem_cache_alloc);
```

```
2916 static inline void *____cache_alloc(struct kmem_cache *cachep, gfp_t flags)
2917 {
2918         void *objp;
2919         struct array_cache *ac;
2920         bool force_refill = false;
2921
2922         check_irq_off();
2923
2924         ac = cpu_cache_get(cachep);
2925         if (likely(ac->avail)) {
2926                 ac->touched = 1;
2927                 objp = ac_get_obj(cachep, ac, flags, false);
2928
2929                 /*
2930                  |* Allow for the possibility all avail objects are not allowed
2931                  |* by the current flags
2932                  |*/
2933                 if (objp) {
2934                         STATS_INC_ALLOCHIT(cachep);
2935                         goto out;
2936                 }
2937                 force_refill = true;
2938         }
2939
2940         STATS_INC_ALLOCMISS(cachep);
2941         objp = cache_alloc_refill(cachep, flags, force_refill);
2942         /*
2943          |* the 'ac' may be updated by cache_alloc_refill(),
2944          |* and kmemleak_erase() requires its correct value.
2945          |*/
2946         ac = cpu_cache_get(cachep);
```

# 例子1：笨老师，为啥我建的小仓库在slabinfo中没找到？

```
/ # cat /proc/slabinfo
slabinfo - version: 2.1
# name            <active_objs> <num_objs> <objsize> <objperslab> <pagesperslab> : tunables <limit> <batchcount> <sharedfactor> : slabdata <a
ctive_slabs> <num_slabs> <sharedavail>
ubifs_inode_slab        0       0     424     19      2 : tunables       0      0      0 : slabdata       0      0      0
v9fs_inode_cache        0       0     368     22      2 : tunables       0      0      0 : slabdata       0      0      0
jffs2_refblock          0       0     248     16      1 : tunables       0      0      0 : slabdata       0      0      0
jffs2_full_dnode        0       0      16    256      1 : tunables       0      0      0 : slabdata       0      0      0
jffs2_i                 0       0     376     21      2 : tunables       0      0      0 : slabdata       0      0      0
nfs_direct_cache        0       0     176     23      1 : tunables       0      0      0 : slabdata       0      0      0
nfs_inode_cache         0       0     584     28      4 : tunables       0      0      0 : slabdata       0      0      0
fat_inode_cache         0       0     416     19      2 : tunables       0      0      0 : slabdata       0      0      0
fat_cache               0       0      24    170      1 : tunables       0      0      0 : slabdata       0      0      0
```

# 使用slabinfo来查看是否被merged？

```
/mnt # slabinfo -a

:at-0000016   <- jbd2_revoke_table_s isp1760_urb_listitem revoke_record revoke_table
:at-0000032   <- jbd2_journal_handle ext4_extent_status jbd2_revoke_record_s isp1760_qh
:at-0000040   <- ext4_io_end isp1760_qtd ext4_free_data
:at-0000064   <- jbd2_journal_head journal_head
:t-0000024    <- scsi_data_buffer jbd2_inode nsproxy jffs2_node_frag ip_fib_alias ubi_wl_entry_slab jffs2_inode_cache dnotify_struct
:t-0000032    <- ip_fib_trie ftrace_event_field sd_ext_cdb figo-cache secpath_cache anon_vma_chain fasync_cache tcp_bind_bucket file_lock_ctx ext4_system_zone
:t-0000040    <- eventpoll_pwq jffs2_tmp_dnode
:t-0000064    <- fs_cache nfs_page blkdev_ioc inotify_inode_mark dnotify_mark pid kmalloc-64 jffs2_raw_dirent uid_cache kiocb
:t-0000088    <- flow_cache vm_area_struct
:t-0000128    <- cred_jar sgpool-8 jffs2_raw_inode rpc_tasks bio_integrity_payload eventpoll_epi inet_peer_cache bio-0 kmalloc-128 file_lock_cache ip_dst_cache
:t-0000192    <- skbuff_head_cache mnt_cache filp request_sock_TCP virtio_scsi_cmd biovec-16 kmalloc-192
:t-0000256    <- sgpool-16 kmalloc-256 pool_workqueue files_cache
:t-0000384    <- skbuff_fclone_cache dio
:t-0000448    <- kioctx nfs_commit_data mm_struct
:t-0000512    <- sgpool-32 kmalloc-512
:t-0000576    <- RAW nfs_read_data signal_cache nfs_write_data PING UNIX
:t-0001024    <- kmalloc-1024 sgpool-64
:t-0002048    <- sgpool-128 kmalloc-2048 rpc_buffers
:t-0004096    <- names_cache kmalloc-4096
/mnt # slabinfo -a | grep figo
:t-0000032    <- ip_fib_trie ftrace_event_field sd_ext_cdb figo-cache secpath_cache anon_vma_chain fasync_cache tcp_bind_bucket file_lock_ctx ext4_system_zone
/mnt #
```

# slabtop工具

# 怀疑slab吃掉很多内存？用nmon看看

# 使用slub_debug进行内存泄漏检测

- 详情见《奔跑吧Linux内核》第6.4.1章内容

# Thanks

奔跑吧Linux社区出品

开源精神 奔跑不息