

# Page Allocation

SWE3015

Jeaho Hwang

- Unit of memory mapping
  - 4KB size by default
  - For every page there is ***struct page***.
    - $36\text{B} / 4\text{KB} = 0.88\%$  overhead

```
struct page {  
    unsigned long flags;  
    struct address_space *mapping;  
    pgoff_t index;  
    atomic_t _mapcount;  
    atomic_t _count;  
    struct list_head lru;  
    unsigned long private;  
    void *virtual;  
}  
  
<include/linux/mm_types.h>
```



- Flag field
  - Stores the status of the page
    - Dirty, locked, etc..
  - Enumeration type defined in include/linux/page-flags.h
- Count field
  - Reference count for the page
  - Accessed with wrapper function
    - Page\_count()

```
static inline int page_count(struct page *page)
{
    return atomic_read(&compound_head(page)-
>_count);
}
<include/linux/mm.h>
```

```
enum pageflags {
    PG_locked,
    PG_error,
    PG_referenced,
    PG_uptodate,
    PG_dirty,
    PG_lru,
    ...
}
<include/linux/page-flags.h>
```



- Virtual field
  - Page's virtual address
- Lru field
  - Managing page list
  - Use in page replacement
- Others...
  - mapping



# Zones

- To group pages of similar properties
- Types
  - ZONE\_DMA, ZONE\_DMA32
    - Can be used for Direct Memory Access
  - ZONE\_NORMAL
    - Regularly mapped pages
  - ZONE\_HIGHMEM
    - Over 896MB for i386
  - ZONE\_MOVABLE

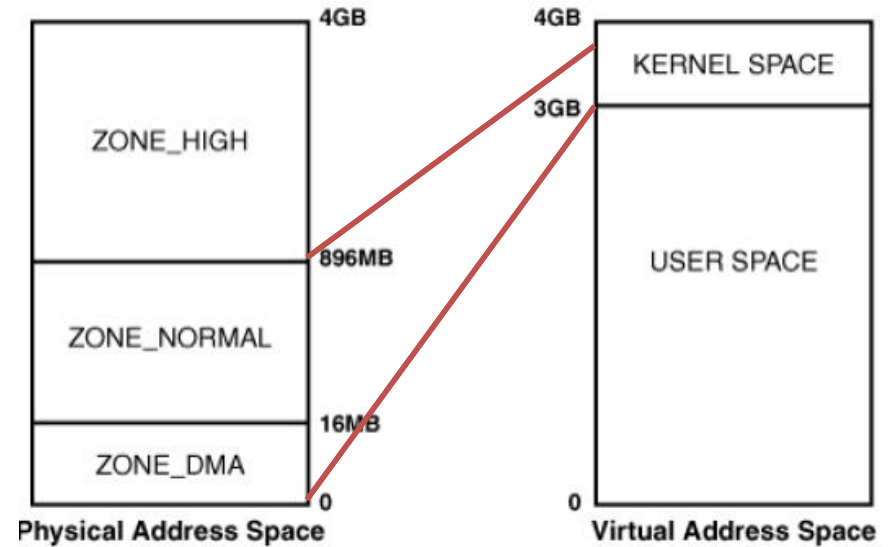
```
Enum zone_type
{
    ZONE_DMA,
    ZONE_DMA32,
    ZONE_NORMAL,
    ZONE_HIGHMEM,
    ZONE_MOVABLE,
    __MAX_NR_ZONES
}

<include/linux/mmzone.h>
```



# Basic structure - Zone(cont'd)

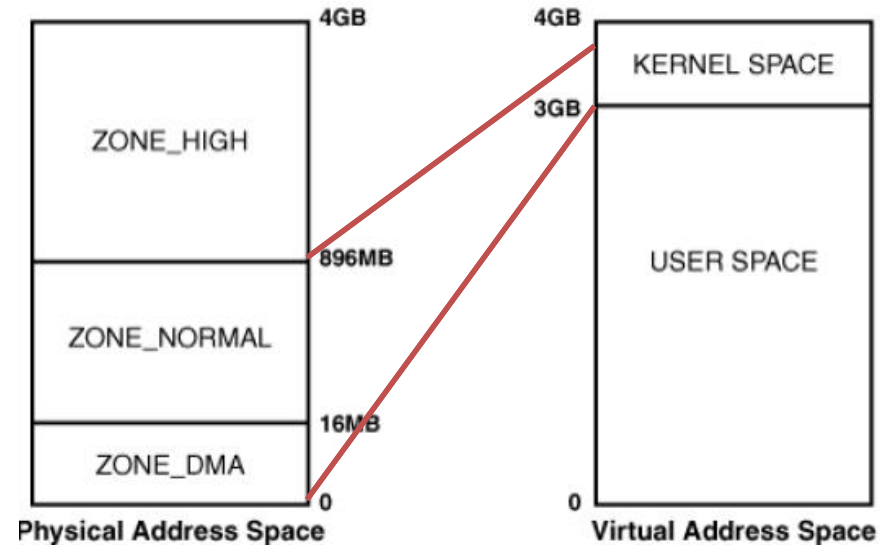
- ZONE\_DMA
  - ISA can DMA only in the first 16MB area
    - For those devices
- ZONE\_DMA\_32
  - For devices only supports 32bit address space in x86\_64
    - So it provides the first 4GB area.





# Basic structure - Zone(cont'd)

- ZONE\_NORMAL
  - Address which can be accessed directly by kernel
    - Except the reservation, 0 ~ 896MB area in i386
    - 0 ~ 64GB in x86\_64



- ZONE\_HIGHMEM
  - Address which cannot be accessed directly by kernel
    - Needs special operation to map in kernel space
    - 896MB ~ 4GB in i386
    - **No HIGHMEM for x86\_64** since kernel can cover all memory directly



# Getting Pages

- Allocating memory with page-size granularity
  - `<linux/gfp.h>`
    - `alloc_pages(gfp_t gfp_mask, unsigned int order)`
      - $2^{\text{order}}$  core allocation function
    - `void *page_address(const struct page *page)`
      - Get address of a page structure
    - `__get_free_pages`
      - Getting pages be `alloc_pages()` and return using `page_address()`
    - `get_zeroed_page(gfp_t gfp_mask)`
      - Getting an empty page
    - `#define alloc_page(gfp_mask) alloc_pages(gfp_mask, 0)`
    - `#define __get_free_page(gfp_mask)`  
`__get_free_pages(gfp_mask, 0)`





# Freeing Pages

- `__free_pages(struct page *page, unsigned int order)`
  - Passing page structure for argument
- `free_pages(unsigned long addr, unsigned int order)`
  - Passing address for argument
- `free_page(unsigned long addr)`



# kmalloc

- `void *kmalloc(size_t size, gfp_t flags)`
  - Similar operation to `malloc()` in user space
  - To allocate physically continuous memory
  - Can allocate upto 4MB continuous memory

```
#define MAX_ORDER 11  
    <include/linux/mmzone.h>
```

```
#define PAGE_SHIFT 12  
    <arch/x86/include/asm/page_types.h>
```

```
#define KMALLOC_SHIFT_HIGH ((MAX_ORDER + PAGE_SHIFT -1 <=  
25?W  
    (MAX_ORDER + PAGE_SHIFT -1) : 25)  
#define KMALLOC_MAX_SIZE (1UL << KMALLOC_SHIFT_HIGH)  
    <include/linux/mmzone.h>
```



# kmalloc

- Flags
  - Represented by the gfp\_t type (`__gget_free_pages`)
  - Three categories
    - Action modifiers
      - **How** the kernel supposed to allocate
    - Zone modifiers
      - **Where** to allocate
    - Types
      - Combination of action & zone modifiers to stand a certain **type**



# GFP Flags: Action Modifiers

Flag	Description
<code>__GFP_WAIT</code>	The allocator can sleep.
<code>__GFP_HIGH</code>	The allocator can access emergency pools.
<code>__GFP_IO</code>	The allocator can start disk I/O.
<code>__GFP_FS</code>	The allocator can start filesystem I/O.
<code>__GFP_COLD</code>	The allocator should use cache cold pages.
<code>__GFP_NOWARN</code>	The allocator does not print failure warnings.
<code>__GFP_REPEAT</code>	The allocator repeats the allocation if it fails, but the allocation can potentially fail.
<code>__GFP_NOFAIL</code>	The allocator indefinitely repeats the allocation. The allocation cannot fail.
<code>__GFP_NORETRY</code>	The allocator never retries if the allocation fails.
<code>__GFP_NOMEMALLOC</code>	The allocator does not fall back on reserves.
<code>__GFP_HARDWALL</code>	The allocator enforces “hardwall” cpuset boundaries.
<code>__GFP_RECLAIMABLE</code>	The allocator marks the pages reclaimable.
<code>__GFP_COMP</code>	The allocator adds compound page metadata (used internally by the <code>budget</code> code)



# GFP Flags: Zone Modifiers

- Flags can be combined
  - `ptr = kmalloc(size, __GFP_WAIT | __GFP_IO | __GFP_FS);`

---

Flag	Description
<code>__GFP_DMA</code>	Allocates only from <code>ZONE_DMA</code>
<code>__GFP_DMA32</code>	Allocates only from <code>ZONE_DMA32</code>
<code>__GFP_HIGHMEM</code>	Allocates from <code>ZONE_HIGHMEM</code> or <code>ZONE_NORMAL</code>

---



# GFP Flags: types

Flag	Description
<code>GFP_ATOMIC</code>	The allocation is high priority and must not sleep. This is the flag to use in interrupt handlers, in bottom halves, while holding a spinlock, and in other situations where you cannot sleep.
<code>GFP_NOWAIT</code>	Like <code>GFP_ATOMIC</code> , except that the call will not fallback on emergency memory pools. This increases the likelihood of the memory allocation failing.
<code>GFP_NOIO</code>	This allocation can block, but must not initiate disk I/O. This is the flag to use in block I/O code when you cannot cause more disk I/O, which might lead to some unpleasant recursion.
<code>GFP_NOFS</code>	This allocation can block and can initiate disk I/O, if it must, but it will not initiate a filesystem operation. This is the flag to use in filesystem code when you cannot start another filesystem operation.
<code>GFP_KERNEL</code>	This is a normal allocation and might block. This is the flag to use in process context code when it is safe to sleep. The kernel will do whatever it has to do to obtain the memory requested by the caller. This flag should be your default choice.
<code>GFP_USER</code>	This is a normal allocation and might block. This flag is used to allocate memory for user-space processes.
<code>GFP_HIGHUSER</code>	This is an allocation from <code>ZONE_HIGHMEM</code> and might block. This flag is used to allocate memory for user-space processes.
<code>GFP_DMA</code>	This is an allocation from <code>ZONE_DMA</code> . Device drivers that need DMA-able memory use this flag, usually in combination with one of the preceding flags.



# GFP Flags: types

- Composition of type flags

Flag	Value
GFP_ATOMIC	__GFP_HIGH
GFP_NOIO	__GFP_WAIT
GFP_NOFS	(__GFP_WAIT   __GFP_IO)
GFP_KERNEL	(__GFP_WAIT   __GFP_IO   __GFP_FS)
GFP_USER	(__GFP_WAIT   __GFP_IO   __GFP_FS)
GFP_DMA	__GFP_DMA



# GFP Flags

- Which flag to use when?

---

Situation	Solution
Process context, can sleep	Use <code>GFP_KERNEL</code> .
Process context, cannot sleep	Use <code>GFP_ATOMIC</code> , or perform your allocations with <code>GFP_KERNEL</code> at an earlier or later point when you can sleep.
Interrupt handler	Use <code>GFP_ATOMIC</code> .
Softirq	Use <code>GFP_ATOMIC</code> .
Tasklet	Use <code>GFP_ATOMIC</code> .
Need DMA-able memory, can sleep	Use <code>(GFP_DMA   GFP_KERNEL)</code> .
Need DMA-able memory, cannot sleep	Use <code>(GFP_DMA   GFP_ATOMIC)</code> , or perform your allocation at an earlier point when you can sleep.

---





# kfree

- void kfree(const void \*ptr)
  - DO NOT kfree memory **not previously allocated by kmalloc**
    - It may cause bugs



# vmalloc/vfree

- `void *vmalloc(unsigned long size)`
  - Allocating more than 4MB
  - Logically continuous, but physically not
  - Large overhead
    - Changing kernel page table to allocate non-continuous area
    - Use `kmalloc` as possible
    - Kernel modules are loaded into memory via `vmalloc`
- `void vfree(const void *addr)`
- Declared in `mm/vmalloc.c`



# Slab Layer

- Why slab layer is required?
  - A lot of data structures are frequently allocated/freed.
    - By naïve allocation, slowdown/fragmentation caused
  - To solve it, *free list* is maintained for each structure.
    - A block of available, already allocated data structures
  - There exists no central control by kernel for free lists.
    - E.g. shrink the list size if available memory size is low
- **The slab layer is a generic data structure-caching layer**
  - task\_struct, inode, mm\_struct, etc.



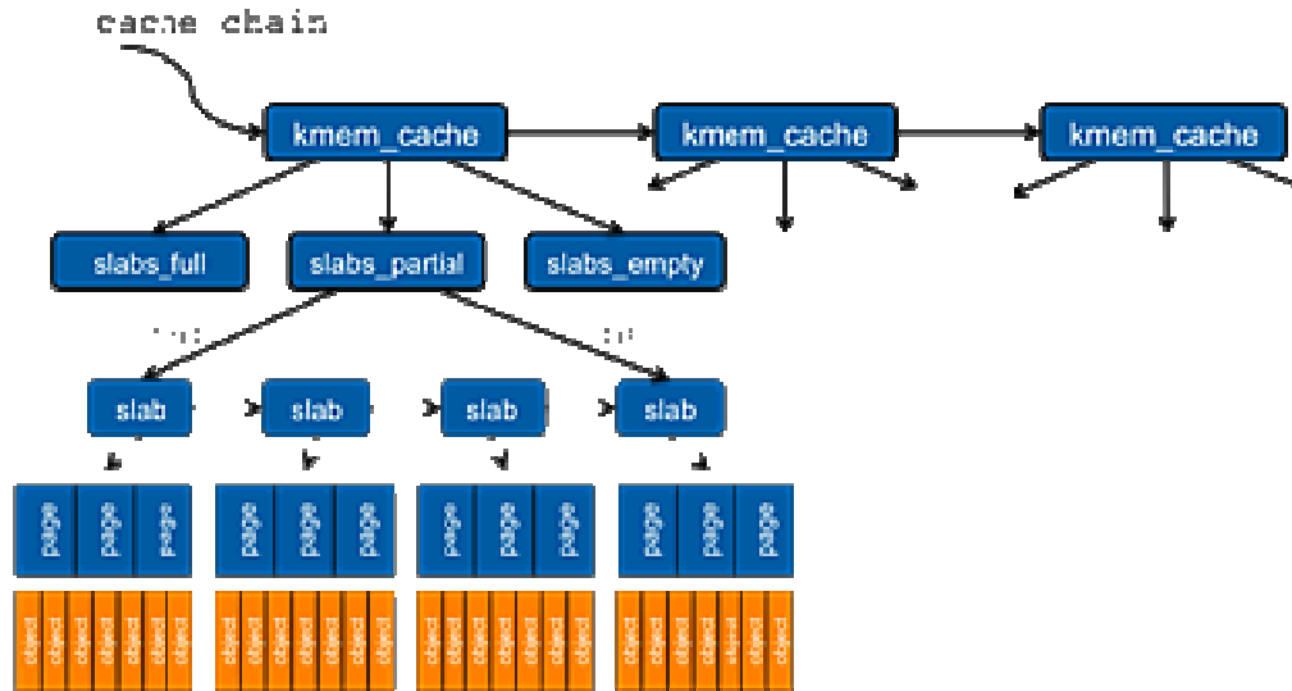
# Slab Layer

- Design of the slab layer
  - *Cache*: a storage for a specific type of object
    - One cache per object type
    - semaphores, file objects, process descriptors, etc.
    - `kmalloc()` is built on the slab layer
  - *Slab*: a contiguous piece of memory, often several page size.
    - A cache is stored in 1 or more slabs.
    - Each slab contains some number of equal-sized *objects*.
      - No fragmentation
    - Three states
      - Full: all objects in the slab are in use.
      - Empty: all objects in the slab are free, so reclaimable by the kernel.
      - Partial: the slab contains both free and in-use objects.



# Slab Layer

- Design of the slab layer (cont'd)
  - A linked list of caches





# Slab Layer

- Slab operations
  - `kmem_cache_create()`
    - Creating a new cache
    - Typically used when the kernel initializes or a kernel module is loaded
  - `kmem_cache_destroy()`
    - Destroying a cache
  - `void * kmem_cache_alloc(struct kmem_cache *cachep, gfp_t flags)`
    - getting a free object pointer from cachep
    - If no free object, it obtains new pages via `kmem_getpages()`.
  - `void kmem_cache_free(struct kmem_cache *cachep, void *objp)`
    - Freeing objp in cachep

- An example of using the slab allocator

```
hahaman5@ubuntu: ~/linux-3.8.0
void __init fork_init(unsigned long mempages)
{
#ifdef CONFIG_ARCH_TASK_STRUCT_ALLOCATOR
#ifdef ARCH_MIN_TASKALIGN
#define ARCH_MIN_TASKALIGN    L1_CACHE_BYTES
#endif
    /* create a slab on which task_structs can be allocated */
    task_struct_cachep =
        kmem_cache_create("task_struct", sizeof(struct task_struct),
            ARCH_MIN_TASKALIGN, SLAB_PANIC | SLAB_NOTRACK, NULL);
#endif

    /* do the arch specific task caches init */
"kernel/fork.c" 1941 lines --13%--
```

```
hahaman5@ubuntu: ~/linux-3.8.0
#ifdef CONFIG_ARCH_TASK_STRUCT_ALLOCATOR
static struct kmem_cache *task_struct_cachep;

static inline struct task_struct *alloc_task_struct_node(int node)
{
    return kmem_cache_alloc_node(task_struct_cachep, GFP_KERNEL, node);
}

static inline void free_task_struct(struct task_struct *tsk)
{
    kmem_cache_free(task_struct_cachep, tsk);
"kernel/fork.c" 1941 lines --6%--
```

- Checking slab

```

hahaman5@ubuntu: ~/linux-3.8.0
hahaman5@ubuntu:~/linux-3.8.0$ sudo cat /proc/slabinfo
slabinfo - version: 2.1
# name          <active_objs> <num_objs> <objsize> <objperslab> <pagesperslab> : tunables <limit> <batchcount>
t> <sharedfactor> : slabdata <active_slabs> <num_slabs> <sharedavail>
ext2_inode_cache  21      21      752    21      4 : tunables  0  0  0 : slabdata  1  1  0
UDPLITEv6         0       0     1088   15      4 : tunables  0  0  0 : slabdata  0  0  0
UDIPv6           15      15     1088   15      4 : tunables  0  0  0 : slabdata  1  1  0
tw_sock_TCPv6    0       0      256    16      1 : tunables  0  0  0 : slabdata  0  0  0
TCPv6            8       8     1984    8      4 : tunables  0  0  0 : slabdata  1  1  0
zcache_objnode   0       0      536    15      2 : tunables  0  0  0 : slabdata  0  0  0
kcopyd_job       0       0     3240   10      8 : tunables  0  0  0 : slabdata  0  0  0
dm_uevent        0       0     2608   12      8 : tunables  0  0  0 : slabdata  0  0  0
dm_rq_clone_bio_info 0       0      144    28      1 : tunables  0  0  0 : slabdata  0  0  0
dm_rq_target_io  0       0      416    19      2 : tunables  0  0  0 : slabdata  0  0  0
bsg_cmd          0       0      312    13      1 : tunables  0  0  0 : slabdata  0  0  0
mqueue_inode_cache 9       9      896    9       2 : tunables  0  0  0 : slabdata  1  1  0
fuse_request     0       0      608    13      2 : tunables  0  0  0 : slabdata  0  0  0
fuse_inode       0       0      704    11      2 : tunables  0  0  0 : slabdata  0  0  0
ecryptfs_key_record_cache 0       0      576    14      2 : tunables  0  0  0 : slabdata  0  0  0
ecryptfs_inode_cache 0       0      960    8       2 : tunables  0  0  0 : slabdata  0  0  0
fat_inode_cache  0       0      688    23      4 : tunables  0  0  0 : slabdata  0  0  0
fat_cache        0       0       40    102     1 : tunables  0  0  0 : slabdata  0  0  0
hugetlbfs_inode_cache 14      14     576    14      2 : tunables  0  0  0 : slabdata  1  1  0
journal_handle   170     170     24    170     1 : tunables  0  0  0 : slabdata  1  1  0
journal_head     936     936    112    36      1 : tunables  0  0  0 : slabdata  26  26  0
ext4_inode_cache 31994   31994   912    17      4 : tunables  0  0  0 : slabdata 1882 1882  0
ext4_free_data   64      64      64     64     1 : tunables  0  0  0 : slabdata  1  1  0
ext4_allocation_context 60      60     136    30      1 : tunables  0  0  0 : slabdata  2  2  0
ext4_io_end      70      70     1128   14      4 : tunables  0  0  0 : slabdata  5  5  0
ext4_io_page     770     2304    16    256     1 : tunables  0  0  0 : slabdata  9  9  0
extent_status    128     128     32    128     1 : tunables  0  0  0 : slabdata  1  1  0
ext3_inode_cache 0       0      776    10      2 : tunables  0  0  0 : slabdata  0  0  0
ext3_xattr       0       0       88     46     1 : tunables  0  0  0 : slabdata  0  0  0

```





# Kernel stack

- Every active thread has a kernel stack
  - Statically allocated 2 contiguous pages
    - Which stores task\_struct (SEE Lecture 1)
    - Kernel thread uses only the kernel stack.
  - Used when syscall or interrupt
    - Interrupt handler uses the interrupted process.
  - 4k kernel stack option is available
    - To reduce kernel memory space
    - Interrupt stack per CPU is provided for interrupt handlers
      - Some legacy handlers overflow 4k stack.



# High memory mapping

- Kernel can directly access to 1G space
  - Accessing to other part needs mapping
    - Permanent mapping
      - kmap/unmap
    - Temporary mapping
      - kmap(kunmap)\_atomic
      - Must not sleep between map and unmap
- Use 896MB ~ 1GB space to mapping



# Percpu allocation

- Maintaining a counter per CPU
  - No need to use global lock
  - Reducing cache invalidation
- Pros
  - Reduced locking requirement
  - Reduced cache invalidation
- Cons
  - Disabling kernel preemption
  - Can't sleep in using percpu data



# Percpu allocation

```
#define alloc_percpu(type) \
    (typeof(type) __percpu *)__alloc_percpu(sizeof(type), __alignof__(type)) \
    <include/linux/percpu.h>
```

```
#define get_cpu_var(var) ({ \
    \
    preempt_disable(); \
    &__get_cpu_var(var); }) \
\
#define put_cpu_var(var) do { \
    (void)&(var); \
    preempt_enable(); \
} while (0) \
\
<include/linux/percpu.h>
```

```
Void *percpu_ptr; \
Unsigned long *foo; \
\
Percpu_ptr = alloc_percpu(unsigned long); \
If(!percpu_ptr) \
    /* error handling code */ \
\
Foo = get_cpu_var(percpu_ptr); \
/* manipulate foo .. */ \
\
Put_cpu_var(percpu_ptr); \
\
<Example code>
```