

Memory Allocation :Day 1

SWE3015

Sung-hun Kim



- Virtual memory
 - All processes have their own isolated virtual memory
 - Processes can view virtual address space only
 - All memory space are sliced into pages
 - Page table stores virtual-to-physical address mapping
- Demand paging
 - OSes provide lazy page allocation/load for effective memory management
 - Page fault handler handles several *faults*
 - Major fault
 - Minor fault



- Unit of memory mapping
 - 4KB size by default (can be extended 2 MiB, 1 GiB according to hardware specification)
 - For every page there is **struct page**.
 - 36B / 4KiB = 0.88% overhead
 - For 4 GiB memory system -> 20 MiB of page structures

```
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(least recently used) field
 - Managing page list
 - Use in page replacement
- Others...
 - mapping
 - Used for managing address space of device file
 - Or mapped to anon_vma



Kernel Functions for Page Structure

- `get_page(struct page *page)`
 - Increment page's `_count`
- `put_page(struct page *page)`
 - Decrement page's `_count`
- `lock_page(struct page *page)`
 - Lock page by using `page->flag` and `wait_queue`
- `unlock_page(struct page *page)`
 - Unlock page by using `page->flag` and `wait_queue`
 - Wake up waiting processes



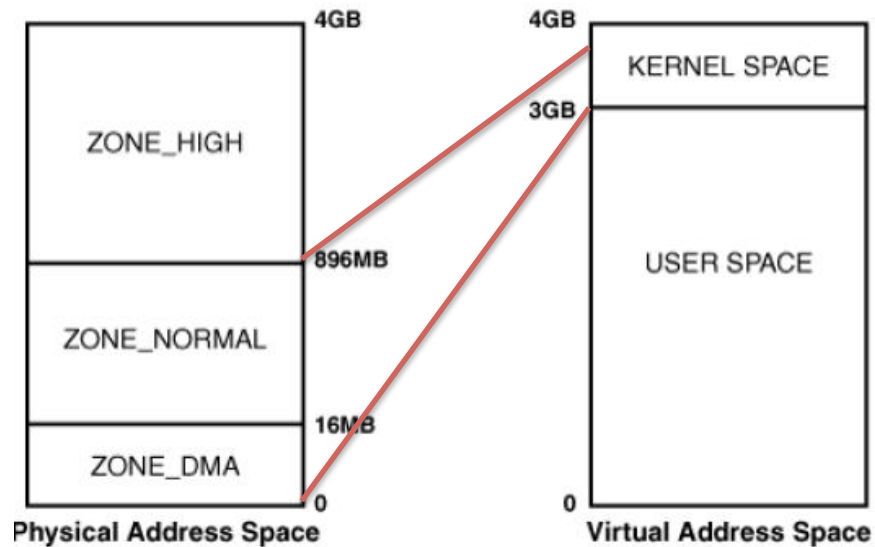
- To group pages of similar properties
- Types
 - ZONE_DMA, ZONE_DMA32
 - Can be used for Direct Memory Access
 - ZONE_NORMAL
 - Regularly mapped pages
 - ZONE_HIGHMEM
 - Over 896MB for i386
 - ZONE_MOVABLE

```
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)

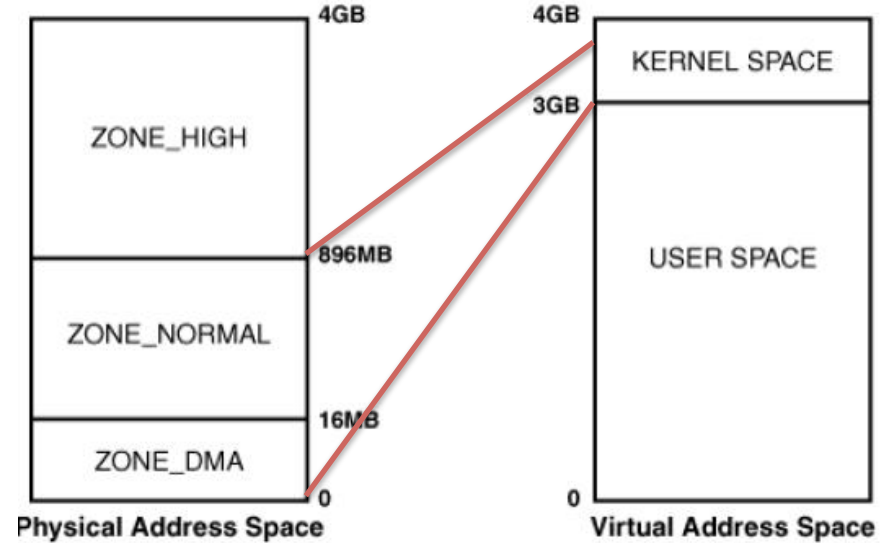
- ZONE_DMA
 - ISA can DMA only in the first 16MB area
 - For those devices
- ZONE_DMA_32
 - For devices only supports 32bit address space in x86_64
 - So it provides the first 4GB area.





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^{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

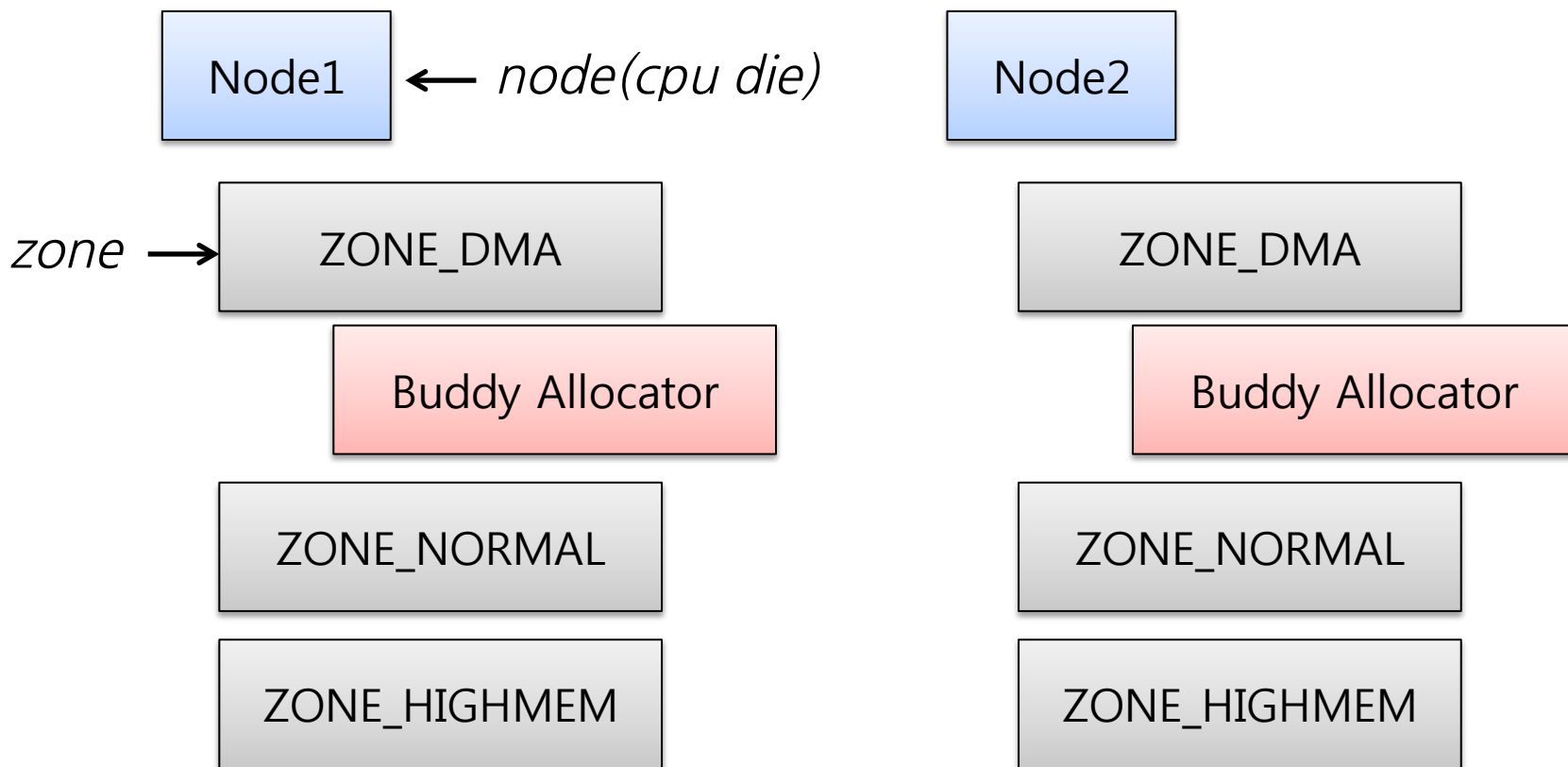
- `__free_page(struct page *page, unsigned int order)`
 - Passing page structure for argument
- `free_page(unsigned long addr, unsigned int order)`
 - Passing address for argument
- `free_page(unsigned long addr)`



Page Allocation Steps

`alloc_pages(gfp, order)`

`__alloc_pages_nodemask(gfp, order, zone, node)`





- 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 <= 2  
5?W  
                                (MAX_ORDER + PAGE_SHIFT -1) : 25)  
#define KMALLOC_MAX_SIZE (1UL << KMALLOC_SHIFT_HIGH)  
                                <include/linux/mmzone.h>
```



- Flags
 - Represented by the gfp_t type (**__**get_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>hugetlb</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



- 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.



- `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`