Linux Kernel Programming
**Memory Management**

Pierre Olivier

Systems Software Research Group @ Virginia Tech

March 28, 2017

Virginia Tech

## Outline

1. Pages and zones

2. Low-level memory allocator

3. `kmalloc()` and `vmalloc()`

4. Slab layer
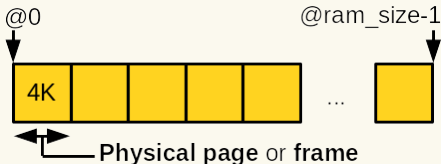
5. Stack, high memory and per-cpu allocation

# Outline

Virginia Tech

# Pages and zones
Pages

- ▶ Memory allocation in the kernel is different from user space
  - ▶ How to handle memory allocation failures?
  - ▶ In some situations the kernel cannot sleep
  - ▶ Etc.
- ▶ The **page** is the basic unit for memory management by the MMU, and thus the kernel
  - ▶ Page size is machine dependent
    - ▶ Typical values for x86 are **4K**, 2M, and 1G
    - ▶ Get the page size on your machine:

```
1  getconf PAGESIZE
```



**Physical page** or **frame**

# Pages and zones
Pages: `struct page`

- ► Each **physical page** is represented by a **`struct page`**
  - ► Most of the pages are used for (1) kernel/userspace memory (*anonymous mapping*) or (2) file mapping

- ► Simplified version:

```
1  struct page {
2    unsigned long flags;
3    unsigned counters;
4    atomic_t _mapcount;
5    unsigned long private;
6    struct address_space *mapping;
7    pgoff_t index;
8    struct list_head lru;
9    void *virtual;
10 }
```

- ► More info on `struct page`: [1]

- ► `flags`: page status (permission, dirty, etc.)
- ► `counters`: usage count
- ► `private`: private mapping
- ► `mapping`: file mapping
- ► `index`: offset within mapping
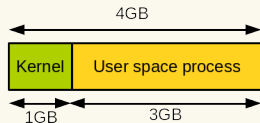- ► `virtual`: virtual address

Virginia Tech

# Pages and zones
## Pages (2)

- The kernel uses `struct page` to keep track of the owner of the page
  - User-space process, kernel statically/dynamically allocated data, page cache, etc.
- **There is one `struct page` object per physical memory page**
  - `sizeof(struct page)` with regular kernel compilation options: 64 bytes
  - Assuming 8GB of RAM and 4K-sized pages: 128MB reserved for `struct page` objects (~1.5%)
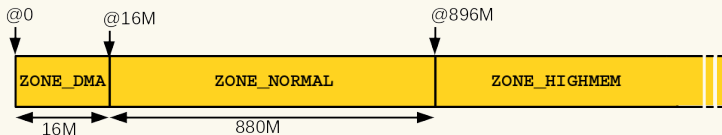
Virginia Tech

# Pages and zones
Zones

- Because of hardware limitations, **only certain physical pages can be used in certain contexts**
    - Example: on x86 some DMA-enabled devices can only access the lowest 16M part of physical memory
    - On x86_32 the kernel address space area sometimes cannot map all physical RAM (1/3 model)
- Physical memory is divided into **zones**:
    - `ZONE_DMA`: pages with which DMA can be used
    - `ZONE_DMA32`: memory for other DMA limited devices
    - `ZONE_NORMAL`: page always mapped in the address space
    - `ZONE_HIGHMEM`: pages only mapped temporary
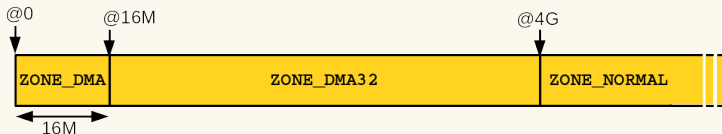- Zones layout is completely architecture dependent

# Pages and zones
## Zones (2)

► x86_32 zones layout:



► x86_64 zones layout:



► Some memory allocation requests must be served in specific zones

► While the kernel tries to avoid it, general allocations requests can be served from any zone if needed

# Pages and zones
## Zones (3)

- ▶ Each zone is represented by a `struct zone` object
  - ▶ Defined in `include/linux/mmzone.h`
  - ▶ Simplified version with notable fields:

```
1  struct zone {
2    unsigned long       watermark[NR_WMARK];
3    const char          *name;
4    spinlock_t          lock;
5    struct free_area    free_area[MAX_ORDER];
6    /* ... */
7  }
```
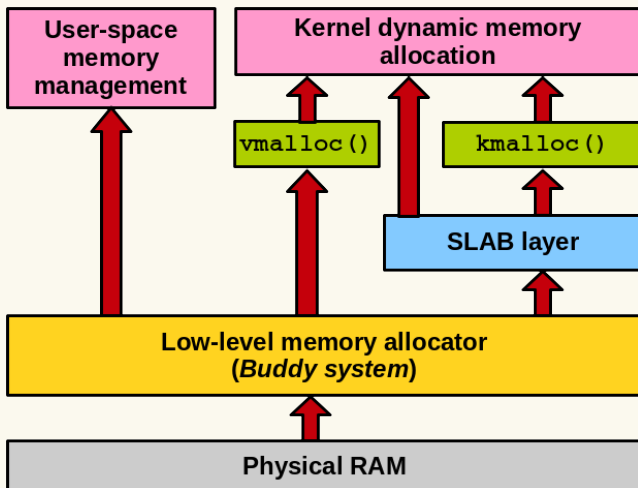
- ▶ `watermark` minimum, low and high watermarks used for per-area memory allocation
- ▶ `lock`: protects against concurrent access
- ▶ `free_area`: list of free pages to serve memory allocation requests
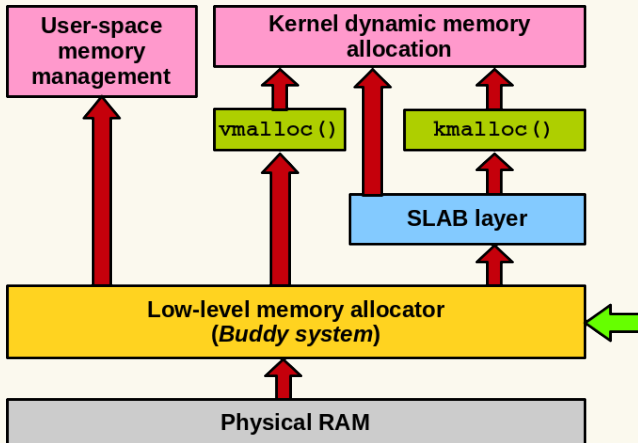
Virginia Tech

# Outline

Virginia
Tech

# Low-level memory allocator

Memory allocation: general overview

# Low-level memory allocator

## Memory allocation: general overview

# Low-level memory allocator
Getting pages

- ▶ Low-level mechanisms allocating memory with page-sized granularity
- ▶ Interface in `include/linux/gfp.h`
  - ▶ Core function:

```
1 struct page * alloc_pages(gfp_t gfp_mask, unsigned int order);
```

  - ▶ Allocates $2^{order}$ contiguous pages ($1 << order$)
  - ▶ Returns the **address of the first allocated `struct page`**
  - ▶ To actually use the allocated memory, need to convert to virtual address:

```
1 void * page_address(struct page *page);
```

  - ▶ To allocate and get the virtual address directly:

```
1 unsigned long __get_free_pages(gfp_t gfp_mask, unsigned int order);
```

# Low-level memory allocator
Getting pages (2)

- ▶ To get a single page:

```
1 struct page * alloc_page(gfp_t gfp_mask);
2 unsigned long __get_free_page(gfp_t gfp_mask);
```

- ▶ To get a page filled with zeros:

```
1 unsigned long get_zeroed_page(gfp_t gfp_mask);
```

  - ▶ Needed for pages given to user space
    - ▶ A page containing user space data (process A) that is freed can be later given to another process (process B)
    - ▶ We don't want process B to read information from process A

Virginia Tech

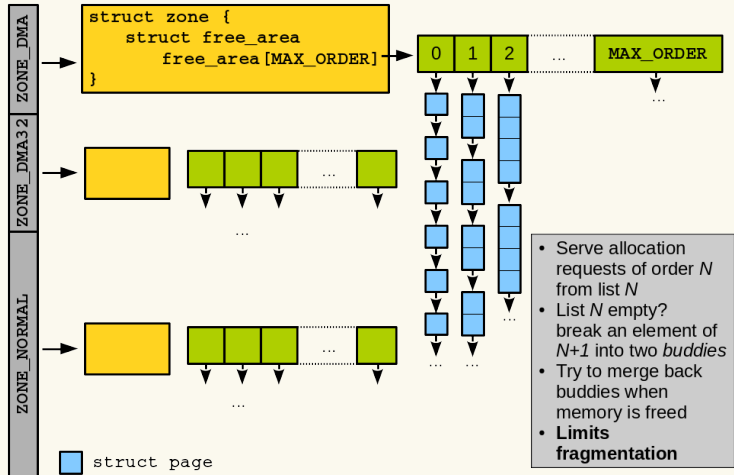# Low-level memory allocator
Freeing pages

▶ Freeing pages - low level API:

```
1  void __free_pages(struct page *page, unsigned int order);
2  void free_pages(unsigned long addr, unsigned int order);
3  void __free_page(struct page *page);
4  void free_page(unsigned long addr);
```

    ▶ Free only the pages you allocate!
        ▶ Otherwise: corruption

Virginia
Tech

# Low-level memory allocator
## Buddy system



struct page

- Serve allocation requests of order *N* from list *N*
- List *N* empty? break an element of *N+1* into two *buddies*
- Try to merge back buddies when memory is freed
- **Limits fragmentation**

# Low-level memory allocator
Usage example

```c
1   #include <linux/module.h>
2   #include <linux/kernel.h>
3   #include <linux/init.h>
4   #include <linux/gfp.h>
5
6   #define PRINT_PREF      "[LOWLEVEL]: "
7   #define PAGES_ORDER_REQUESTED 3
8   #define INTS_IN_PAGE (PAGE_SIZE/sizeof(int))
9
10  unsigned long virt_addr;
11
12  static int __init my_mod_init(void)
13  {
14    int *int_array;
15    int i;
16
17    printk(PRINT_PREF "Entering module.\n");
18
19    virt_addr = __get_free_pages(GFP_KERNEL,
        PAGES_ORDER_REQUESTED);
20    if(!virt_addr) {
21      printk(PRINT_PREF "Error in allocation\n
        ");
22      return -1;
23    }
```

```c
24
25    int_array = (int *)virt_addr;
26    for(i=0; i<INTS_IN_PAGE; i++)
27      int_array[i] = i;
28
29    for(i=0; i<INTS_IN_PAGE; i++)
30      printk(PRINT_PREF "array[%d] = %d\n", i,
        int_array[i]);
31
32    return 0;
33  }
34
35  static void __exit my_mod_exit(void)
36  {
37    free_pages(virt_addr,
        PAGES_ORDER_REQUESTED);
38    printk(PRINT_PREF "Exiting module.\n");
39  }
40
41  module_init(my_mod_init);
42  module_exit(my_mod_exit);
```

Virginia Tech

# Outline

Virginia Tech

# kmalloc() and vmalloc()
kmalloc() usage and kfree()

- ► kmalloc() allocates byte-sized chunks of memory
- ► Allocated memory is **physically contiguous**
- ► Usage similar to userspace malloc()
  - ► Returns a pointer to the first allocated byte on success
  - ► Returns NULL on error
- ► Declared in includes/linux/slab.h:

```
1  void * kmalloc(size_t size, gfp_t flags);
```

- ► Usage example:

```
1  struct my_struct *ptr;
2
3  ptr = kmalloc(sizeof(struct my_struct), GFP_KERNEL);
4  if(!ptr) {
5    /* handle error */
6  }
```

Virginia Tech

# kmalloc() and vmalloc()
## Kmalloc flags

- ▶ gfp_t type (unsigned int in include/linux/types.h)
    - ▶ *get free page* [2]
- ▶ Specify options for the allocated memory
    1. **Action modifiers**
        - ▶ How should the memory be allocated? (for example, can the kernel sleep during allocation?)
    2. **Zone modifiers**
        - ▶ From which zone should the allocated memory come
    3. **Type flags**
        - ▶ Combination of action and zone modifiers
        - ▶ Generally preferred compared to direct use of action/zone

Virginia Tech

# kmalloc() and vmalloc()

Kmalloc flags: action modifiers

- **Action modifiers:**

| Flag | Description |
|------|-------------|
| __GFP_WAIT | Allocator can sleep |
| __GFP_HIGH | Allocator can access emergency pools |
| __GFP_IO | Allocator can start disk IO |
| __GFP_FS | Allocator can start filesystem IO |
| __GFP_NOWARN | Allocator does not print failure warnings |
| __GFP_REPEAT | Allocator repeats the allocation if it fails, the allocation can potentially fail |
| __GFP_NOFAIL | Allocator indefinitely repeats the allocation, which cannot fail |

# kmalloc() and vmalloc()
Kmalloc flags: action modifiers (2)

| Flag | Description |
|------|-------------|
| __GFP_NORETRY | Allocator does not retry if the allocation fails |
| __GFP_NOMEMALLOC | Allocation does not fall back on reserves |
| __GFP_HARDWALL | Allocator enforces "hardwall" cpuset boundaries |
| __GFP_RECLAIMABLE | Allocator marks the pages reclaimable |
| __GFP_COMP | Allocator adds compound page metadata (used by hugetlb code) |

▶ Several action modifiers can be specified together:

```
1  ptr = kmalloc(size, __GFP_WAIT | __GFP_FS | __GFP_IO);
```

    ▶ Generally, **type modifiers are used instead**

Virginia Tech

# kmalloc() and vmalloc()

Kmalloc flags: zone modifiers

- **Zone modifiers:**

| Flag | Description |
|------|-------------|
| __GFP_DMA | Allocates **only** from ZONE_DMA |
| __GFP_DMA32 | Allocates **only** from ZONE_DMA32 |
| __GFP_HIGHMEM | Allocated from ZONE_HIGHMEM or ZONE_NORMAL |

- __GFP_HIGHMEM: OK to use high memory, normal works too
- No flag specified?
    - Kernel allocates from ZONE_NORMAL or ZONE_DMA with a strong preference for ZONE_NORMAL
- Cannot specify __GFP_HIGHMEM to kmalloc() or __get_free_pages()
    - These function return a virtual address
    - Addresses in ZONE_HIGHMEM might not be mapped yet

# kmalloc() and vmalloc()
Kmalloc flags: type flags

- **GFP_ATOMIC**: high priority allocation that cannot sleep
  - Use in interrupt context, while holding locks, etc.
  - Modifier flags: (__GFP_HIGH)
- **GFP_NOWAIT**: same as GFP_ATOMIC but does not fall back on emergency memory pools
  - More likely to fail
  - Modifier flags: (0)
- **GFP_NOIO**: can block but does not initiate disk IO
  - Used in block layer code to avoid recursion
  - Modifier flags: (__GFP_WAIT)
- **GFP_NOFS**: can block and perform disk IO, but does not initiate filesystem operations
  - Used in filesystem code
  - Modifier flags: (__GFP_WAIT | __GFP_IO)

Virginia Tech

# kmalloc() and vmalloc()
Kmalloc flags: type flags (2)

- **GFP_KERNEL**: default choice, can sleep and perform IO
  - Modifier flags: (__GFP_WAIT | __GFP_IO | __GFP_FS)
- **GFP_USER**: normal allocation, might block, for user-space memory
  - Modifier flags: (__GFP_WAIT | __GFP_IO | __GFP_FS)
- **GFP_HIGHUSER**: allocation for user-space memory, from ZONE_HIGHMEM
  - Modifier flags: (__GFP_WAIT | __GFP_IO | __GFP_FS | __GFP_HIGHMEM)
- **GFP_DMA**: served from ZONE_DMA
  - Modifier flags: (__GFP_DMA)

# kmalloc() and vmalloc()
Kmalloc flags: which flag to use?

| Context | Solution |
|---|---|
| Process context, can sleep | GFP_KERNEL |
| Process context, cannot sleep | GFP_ATOMIC or allocate at a different time |
| Interrupt handler | GFP_ATOMIC |
| Softirq | GFP_ATOMIC |
| Tasklet | GFP_ATOMIC |
| Need to handle DMA, can sleep | (GFP_DMA \| GFP_KERNEL) |
| Need DMA, cannot sleep | (GFP_DMA \| GFP_ATOMIC) |

▶ Other types and modifier are declared and documented in
  include/linux/gfp.h

# kmalloc() and vmalloc()
kfree

- ▶ Memory allocated with kmalloc() needs to be freed with kfree()

- ▶ include/linux/slab.h:

```
1 void kfree(const void *ptr);
```

- ▶ Example:

```
1  struct my_struct ptr;
2
3  ptr = kmalloc(sizeof(struct my_struct));
4  if(!ptr) {
5    /* handle error */
6  }
7
8  /* work with ptr */
9
10 kfree(ptr);
```

- ▶ **Do not free memory that has already been freed**
  - ▶ Leads to corruption

Virginia Tech

# kmalloc() and vmalloc()
vmalloc()

- ▶ **vmalloc() allocates *virtually* contiguous pages that are not guarantee to map to *physically* contiguous ones**
    - ▶ Page table is modified → no problems from the programmer usability standpoint
- ▶ Buffers related to communication with the hardware generally need to be physically contiguous
    - ▶ But most of the rest do not, for example user-space buffers
- ▶ However, most of the kernel uses kmalloc() for **performance reasons**
    - ▶ Pages allocated with kmalloc() are directly mapped
    - ▶ Less overhead in virtual to physical mapping setup and translation
- ▶ vmalloc() is still needed to allocate large portions of memory

# kmalloc() and vmalloc()
vmalloc() (2)

- **vmalloc() usage:**
- Similar to user-space malloc()
    - include/linux/vmalloc.h:

```
1  void *vmalloc(unsigned long size);
2  void vfree(const void *addr);
```

- Example:

```
1  struct my_struct *ptr;
2
3  ptr = vmalloc(sizeof(struct my_struct));
4  if(!ptr) {
5    /* handle error */
6  }
7
8  /* work with ptr */
9
10 vfree(ptr);
```

Virginia Tech

# kmalloc() and vmalloc()

vmalloc(): kmalloc() allocated size limitation

```c
1  #include <linux/module.h>
2  #include <linux/kernel.h>
3  #include <linux/init.h>
4  #include <linux/slab.h>
5
6  #define PRINT_PREF "[KMALLOC_TEST]: "
7
8  static int __init my_mod_init(void)
9  {
10    unsigned long i;
11    void *ptr;
12
13    printk(PRINT_PREF "Entering module.\n");
14
15    for(i=1;;i*=2) {
16      ptr = kmalloc(i, GFP_KERNEL);
17      if(!ptr) {
18        printk(PRINT_PREF "could not
         allocate %lu bytes\n", i);
19        break;
20      }
21      kfree(ptr);
22    }
23
24    return 0;
25  }
```

```c
26  static void __exit my_mod_exit(void)
27  {
28    printk(KERN_INFO "Exiting module.\n");
29  }
30
31  module_init(my_mod_init);
32  module_exit(my_mod_exit);
33
34  MODULE_LICENSE("GPL");
```

Virginia Tech

# kmalloc() and vmalloc()

vmalloc(): kmalloc() allocated size limitation (2)



- Max size: 8MB

# Outline

Virginia
Tech

# Slab layer
Presentation

- ▶ Allocating/freeing data structures is done very often in the kernel
- ▶ **Free lists:**
  - ▶ Block of pre-allocated memory for a given type of data structure



- ▶ **Allocation cache**
  - ▶ "Allocating" from the free list is just picking a free slot
  - ▶ "Freeing" from the free list is setting back the slot state to free
- ▶ Faster than frequent $\{k|v\}$malloc() and $\{k|v\}$free()
- ▶ Issue with ad-hoc free lists: no global control

# Slab layer
Presentation (2)

- ▶ **Slab layer/slab allocator**: Linux's generic allocation caching interface
- ▶ Basic principles:
    - ▶ Frequently used data structures are allocated and freed often → needs to be cached
    - ▶ Frequent allocation/free generates memory fragmentation
        - ▶ Caching allocation/free operations in large chunks of contiguous memory reduces fragmentation
    - ▶ An allocator aware of data structure size, page size, and total cache size is more efficient
    - ▶ Part of the cache can be made per-cpu to operate without a lock
    - ▶ Allocator should be NUMA-aware and implement cache-coloring

# Slab layer
Presentation (3)

- The slab layer introduces defines the concept of **caches**
  - One cache per type of data structure supported
    - Example: one cache for struct task_struct, one cache for struct inode, etc.
- Each cache contains one or several **slabs**
  - One or several physically contiguous pages
- Slabs contain **objects**
  - The actual data structure slots
- A slab can be in one of three states:
  - Full: all slots used
  - Partial: some, not all slots used
  - Empty: all slots free
- Allocation requests are served from partial slabs if present, or empty slabs → **reduces fragmentation**
  - A new empty slab is actually allocated in case the cache is full

# Slab layer
## Presentation (4)

# Slab layer
Usage

- A new cache is created using:

```
1  struct kmem_cache *kmem_cache_create(const char *name,
2                                       size_t size,
3                                       size_t align,
4                                       unsigned long flags,
5                                       void (*ctor)(void *));
```

- `name`: cache name
- `size`: data structure size
- `align`: offset of the first element within the page (can put 0)

- `ctor`: constructor called when a new data structure is allocated
    - Rarely used, can put `NULL`
- `flags`: settings controlling the cache behavior

Virginia Tech

# Slab layer
Usage: flags

- `SLAB_HW_CACHEALIGN`
    - Each object in a slab is aligned to a cache line
        - Prevent *false sharing*:

# Slab layer
Usage: flags

- `SLAB_HW_CACHEALIGN`
    - Each object in a slab is aligned to a cache line
        - Prevent *false sharing*:



- Increased memory footprint

# Slab layer
## Usage: flags (2)

- `SLAB_POISON`
    - Fill the slab with know values (`a5a5a5a5`) to detect accesses to uninitialized memory
- `SLAB_RED_ZONE`
    - Extra padding around objects to detect buffer overflows
- `SLAB_PANIC`
    - Slab layer panics if the allocation fails
- `SLAB_CACHE_DMA`
    - Allocation made from DMA-enabled memory

# Slab layer
Usage: allocating from the cache, cache destruction

- **Allocating from the cache:**

```
1 void *kmem_cache_alloc(struct kmem_cache *cachep, gfp_t flags);
```

- **To free an object:**

```
1 void kmem_cache_free(struct kmem_cache *cachep, void *objp);
```

- **To destroy a cache:**

```
1 int kmem_cache_destroy(struct kmem_cache *cachep);
```

- kmem_cache_destroy() should only be called when:
    - All slabs in the cache are empty
    - Nobody is accessing the cache during the execution of
      kmem_cache_destroy()
        - Implement synchronization

Virginia Tech

# Slab layer
Usage example

```
1  #include <linux/module.h>
2  #include <linux/kernel.h>
3  #include <linux/init.h>
4  #include <linux/slab.h>
5
6  #define PRINT_PREF  "[SLAB_TEST] "
7
8  struct my_struct {
9    int int_param;
10   long long_param;
11 };
12
13 static int __init my_mod_init(void)
14 {
15   int ret = 0;
16   struct my_struct *ptr1, *ptr2;
17   struct kmem_cache *my_cache;
18
19   printk(PRINT_PREF "Entering module.\n");
20
21   my_cache = kmem_cache_create("pierre-
       cache", sizeof(struct my_struct),
22     0, 0, NULL);
23   if(!my_cache)
24     return -1;
```

```
25   ptr1 = kmem_cache_alloc(my_cache,
        GFP_KERNEL);
26   if(!ptr1){
27     ret = -ENOMEM;
28     goto destroy_cache;
29   }
30
31   ptr2 = kmem_cache_alloc(my_cache,
        GFP_KERNEL);
32   if(!ptr2){
33     ret = -ENOMEM;
34     goto freeptr1;
35   }
36
37   ptr1->int_param = 42;
38   ptr1->long_param = 42;
39   ptr2->int_param = 43;
40   ptr2->long_param = 43;
```

Virginia Tech

# Slab layer
## Usage example (2)

```
41   printk(PRINT_PREF "ptr1 = {%d, %ld} ; ptr2 = {%d, %ld}\n", ptr1->int_param,
42     ptr1->long_param, ptr2->int_param, ptr2->long_param);
43
44   kmem_cache_free(my_cache, ptr2);
45
46 freeptr1:
47   kmem_cache_free(my_cache, ptr1);
48
49 destroy_cache:
50   kmem_cache_destroy(my_cache);
51
52   return ret;
53 }
54
55 static void __exit my_mod_exit(void)
56 {
57   printk(KERN_INFO "Exiting module.\n");
58 }
59
60 module_init(my_mod_init);
61 module_exit(my_mod_exit);
62
63 MODULE_LICENSE("GPL");
```

Virginia
Tech

# Slab layer
Slab allocator variants and additional info

- **Slab allocator variants**:
    - **SLOB** (Simple List Of Blocks):
        - Used in early Linux version (from 1991)
        - Low memory footprint $\rightarrow$ used in embedded systems
    - **SLAB**:
        - Integrated in 1999
        - Cache-friendly
    - **SLUB**:
        - Integrated in 2008
        - Scalability improvements (amongst other things) over SLAB on many-cores
- More info: [5, 3, 4]

# Outline

Virginia
Tech

# Stack, high memory and per-cpu allocation
Allocating on the stack

- ▶ Each process has:
    - ▶ A user-space stack for execution in user space
    - ▶ A kernel stack for execution in the kernel
- ▶ User-space stack is large and grows dynamically
- ▶ **Kernel stack is small and has a fixed-size**
    - ▶ Generally 8KB on 32-bit architectures and 16KB on 64-bit
- ▶ What about interrupt handlers?
    - ▶ Historically using the kernel stack of the interrupted process
    - ▶ Now using a per-cpu stack (1 single page) dedicated to interrupt handlers
    - ▶ Same thing for softirqs
- ▶ **Reduce kernel stack usage to a minimum**
    - ▶ Local variables & function parameters

Virginia Tech

# Stack, high memory and per-cpu allocation
High memory allocation

- ► On x86_32, physical memory above 896MB is not permanently mapped within the kernel address space
    - ► Because of the limited size of the address space and the 1G/3G kernel/user-space physical memory split
- ► **Before usage, pages from highmem must be mapped after allocation**
    - ► Recall that allocation is done through `alloc_pages()` with the `GFP_HIGHMEM` flag
- ► **Permanent mapping (`include/linux/highmem.h`):**

```
1 void *kmap(struct page *page);
```

  - ► Works on low and high memory
  - ► Maps (update the page table) and return the given
  - ► **May sleep**, use only in process context
  - ► Number of permanent mappings is limited, unmap when done:

```
1 void kunmap(struct page *page);
```

# Stack, high memory and per-cpu allocation
## High memory allocation: usage example

```c
1  #include <linux/module.h>
2  #include <linux/kernel.h>
3  #include <linux/init.h>
4  #include <linux/gfp.h>
5  #include <linux/highmem.h>
6
7  #define PRINT_PREF    "[HIGHMEM]: "
8  #define INTS_IN_PAGE  (PAGE_SIZE/sizeof(
        int))
9
10 static int __init my_mod_init(void)
11 {
12   struct page *my_page;
13   void *my_ptr;
14   int i, *int_array;
15
16   printk(PRINT_PREF "Entering module.\n");
17
18   my_page = alloc_page(GFP_HIGHUSER);
19   if(!my_page)
20     return -1;
21
22   my_ptr = kmap(my_page);
23   int_array = (int *)my_ptr;
```

```c
24   for(i=0; i<INTS_IN_PAGE; i++) {
25     int_array[i] = i;
26     printk(PRINT_PREF "array[%d] = %d\n", i,
       int_array[i]);
27   }
28
29   kunmap(my_page);
30   __free_pages(my_page, 0);
31
32   return 0;
33 }
34
35 static void __exit my_mod_exit(void)
36 {
37   printk(PRINT_PREF "Exiting module.\n");
38 }
39
40 module_init(my_mod_init);
41 module_exit(my_mod_exit);
```

Virginia Tech

# Stack, high memory and per-cpu allocation
High memory allocation: temporary mappings

- ▶ **Temporary mappings**
    - ▶ Also called *atomic mappings* as they can be used from interrupt context
    - ▶ Uses a per-cpu pre-reserved slot

```
1 void *kmap_atomic(struct page *page);
```

    - ▶ Disables kernel preemption
    - ▶ Unmap with:

```
1 void kunmap_atomic(void *addr);
```

- ▶ **Do not sleep while holding a temporary mapping**
- ▶ After kunmap_atomic(), the next temporary mapping will overwrite the slot

Virginia Tech

# Stack, high memory and per-cpu allocation
High memory allocation: temporary mappings usage example

```c
1  #include <linux/module.h>
2  #include <linux/kernel.h>
3  #include <linux/init.h>
4  #include <linux/gfp.h>
5  #include <linux/highmem.h>
6
7  #define PRINT_PREF     "[HIGHMEM_ATOMIC]: "
8  #define INTS_IN_PAGE   (PAGE_SIZE/sizeof(
       int))
9
10 static int __init my_mod_init(void)
11 {
12   struct page *my_page;
13   void *my_ptr;
14   int i, *int_array;
15
16   printk(PRINT_PREF "Entering module.\n");
17
18   my_page = alloc_page(GFP_HIGHUSER);
19   if(!my_page)
20     return -1;
21
22   my_ptr = kmap_atomic(my_page);
23   int_array = (int *)my_ptr;
```

```c
24   for(i=0; i<INTS_IN_PAGE; i++) {
25     int_array[i] = i;
26     printk(PRINT_PREF "array[%d] = %d\n", i,
       int_array[i]);
27   }
28
29   kunmap_atomic(my_ptr);
30   __free_pages(my_page, 0);
31
32   return 0;
33 }
34
35 static void __exit my_mod_exit(void)
36 {
37   printk(PRINT_PREF "Exiting module.\n");
38 }
39
40 module_init(my_mod_init);
41 module_exit(my_mod_exit);
```

Virginia Tech

# Stack, high memory and per-cpu allocation
Per-CPU allocation

- ▶ **Per-cpu data: data that is unique to each CPU (i.e. each core)**
    - ▶ No locking required
    - ▶ Implemented through arrays in which each index corresponds to a CPU:

```
1  unsigned long my_percpu[NR_CPUS]; /* NR_CPUS contains the number of cores */
```

```
1  int cpu;
2
3  cpu = get_cpu();    /* get current CPU, disable kernel preemption */
4  my_percpu[cpu]++;   /* access the data */
5  put_cpu();          /* re-enable kernel preemption */
```

- ▶ Disabling kernel preemption (get_cpu()/put_cpu()) while accessing per-cpu data is necessary:
    - ▶ Preemption then reschedule on another core → cpu not valid anymore
    - ▶ Another task preempting the current one might access the per-cpu data → race condition

# Stack, high memory and per-cpu allocation
Per-CPU allocation: the *percpu* API

- ▶ Linux provides an API to manipulate per-cpu data: *percpu*
  - ▶ In `include/linux/percpu.h`
- ▶ **Compile-time per-cpu data structure usage**:
  - ▶ Creation:

```
1  DEFINE_PER_CPU(type, name);
```

  - ▶ To refer to a per-cpu data structure declared elsewhere:

```
1  DECLARE_PER_CPU(name, type);
```

  - ▶ Data manipulation:

```
1  get_cpu_var(name)++;  /* increment 'name' on this CPU */
2  put_cpu_var(name);    /* Done, disable kernel preemption */
```

  - ▶ Access another CPU data:

```
1  per_cpu(name, cpu)++;  /* increment name on the given CPU */
```

  - ▶ **Need to use locking!**

# Stack, high memory and per-cpu allocation
Per-CPU allocation (2)

- **Per-cpu data at runtime**:
    - Allocation:

```
1  struct my_struct *my_var = alloc_percpu(struct my_struct);
2  if(!my_var) {
3      /* allocation error */
4  }
```

    - Manipulation:

```
1  get_cpu_var(my_var)++;
2  put_cpu_var(my_var);
```

    - Deallocation:

```
1  free_percpu(my_var);
```

- Benefits of per-cpu data:
    - Removes/minimizes the need for locking
    - Reduces cache thrashing
        - Processor access local data so there is less cache coherency overhead (invalidation) in multicore systems

# Stack, high memory and per-cpu allocation
## Per-CPU allocation: usage example (static)

```c
1  #include <linux/module.h>
2  #include <linux/kernel.h>
3  #include <linux/init.h>
4  #include <linux/percpu.h>
5  #include <linux/kthread.h>
6  #include <linux/sched.h>
7  #include <linux/delay.h>
8  #include <linux/smp.h>
9
10 #define PRINT_PREF "[PERCPU] "
11 struct task_struct *thread1, *thread2, *
       thread3;
12 DEFINE_PER_CPU(int, my_var);
13
14 static int thread_function(void *data)
15 {
16   while(!kthread_should_stop()) {
17     int cpu;
18     get_cpu_var(my_var)++;
19     cpu = smp_processor_id();
20     printk("cpu[%d] = %d\n", cpu,
         get_cpu_var(my_var));
21     put_cpu_var(my_var);
22     msleep(500);
23   }
24   do_exit(0);
25 }
```

```c
25 static int __init my_mod_init(void)
26 {
27   int cpu;
28
29   printk(PRINT_PREF "Entering module.\n");
30
31   for(cpu=0; cpu<NR_CPUS; cpu++)
32     per_cpu(my_var, cpu) = 0;
33
34   wmb();
35
36   thread1 = kthread_run(thread_function,
       NULL, "percpu-thread1");
37   thread2 = kthread_run(thread_function,
       NULL, "percpu-thread2");
38   thread3 = kthread_run(thread_function,
       NULL, "percpu-thread3");
39
40   return 0;
41 }
```

Virginia Tech

# Stack, high memory and per-cpu allocation
Per-CPU allocation: usage example (static) (2)

```
40  static void __exit my_mod_exit(void)
41  {
42    kthread_stop(thread1);
43    kthread_stop(thread2);
44    kthread_stop(thread3);
45    printk(KERN_INFO "Exiting module.\n");
46  }
47
48  module_init(my_mod_init);
49  module_exit(my_mod_exit);
50
51  MODULE_LICENSE("GPL");
```

Virginia
Tech

# Stack, high memory and per-cpu allocation
Per-CPU allocation: usage example (dynamic)

```c
1  #include <linux/module.h>
2  #include <linux/kernel.h>
3  #include <linux/init.h>
4  #include <linux/percpu.h>
5  #include <linux/kthread.h>
6  #include <linux/sched.h>
7  #include <linux/delay.h>
8  #include <linux/smp.h>
9
10 #define PRINT_PREF "[PERCPU] "
11 struct task_struct *thread1, *thread2, *
       thread3;
12 void *my_var2;
13
14 static int thread_function(void *data)
15 {
16   while(!kthread_should_stop()) {
17     int *local_ptr, cpu;
18     local_ptr = get_cpu_ptr(my_var2);
19     cpu = smp_processor_id();
20     (*local_ptr)++;
21     printk("cpu[%d] = %d\n", cpu, *
       local_ptr);
22     put_cpu_ptr(my_var2);
23     msleep(500);
24   }
25   do_exit(0);
26 }
```

```c
27 static int __init my_mod_init(void)
28 {
29   int *local_ptr;
30   int cpu;
31   printk(PRINT_PREF "Entering module.\n");
32
33   my_var2 = alloc_percpu(int);
34   if(!my_var2)
35     return -1;
36
37   for(cpu=0; cpu<NR_CPUS; cpu++) {
38     local_ptr = per_cpu_ptr(my_var2, cpu);
39     *local_ptr = 0;
40     put_cpu();
41   }
42
43   wmb();
44
45   thread1 = kthread_run(thread_function,
       NULL, "percpu-thread1");
46   thread2 = kthread_run(thread_function,
       NULL, "percpu-thread2");
47   thread3 = kthread_run(thread_function,
       NULL, "percpu-thread3");
48
49   return 0;
50 }
```

# Stack, high memory and per-cpu allocation
Per-CPU allocation: usage example (dynamic) (2)

```
49  static void __exit my_mod_exit(void)
50  {
51    kthread_stop(thread1);
52    kthread_stop(thread2);
53    kthread_stop(thread3);
54
55    free_percpu(my_var2);
56
57    printk(KERN_INFO "Exiting module.\n");
58  }
59
60  module_init(my_mod_init);
61  module_exit(my_mod_exit);
62
63  MODULE_LICENSE("GPL");
```

Virginia Tech

# Stack, high memory and per-cpu allocation
Choosing the right allocation method

- ▶ Need physically contiguous memory?
    - ▶ `kmalloc()` or low-level allocator, with flags:
        - ▶ `GFP_KERNEL` if sleeping is allowed
        - ▶ `GFP_ATOMIC` otherwise
- ▶ Need large amount of memory, not physically contiguous:
    - ▶ `vmalloc()`
- ▶ Frequently creating/destroying large amount of the same data structure:
    - ▶ Use the slab layer
- ▶ Need to allocate from high memory?
    - ▶ Use `alloc_page()` then `kmap()` or `kmap_atomic()`

Virginia Tech

# Bibliography I

[1] Cramming more into struct page.
https://lwn.net/Articles/565097/.
Accessed: 2017-03-12.

[2] Gfp_kernel or slab_kernel?
https://lwn.net/Articles/23042/.
Accessed: 2017-03-18.

[3] The slub allocator.
https://lwn.net/Articles/229984/.
Accessed: 2017-03-18.

[4] BONWICK, J., ET AL.
The slab allocator: An object-caching kernel memory allocator.
In *USENIX summer* (1994), vol. 16, Boston, MA, USA.

[5] LAMETER, C.
Slab allocators in the linux kernel: Slab, slob, slub.
LinuxCon/Dsseldorf 2014,
http://events.linuxfoundation.org/sites/events/files/slides/slaballocators.pdf.
Accessed: 2017-03-18.