# Advanced Operating Systems

# #4

Shinpei Kato

Associate Professor

Department of Computer Science
Graduate School of Information Science and Technology
The University of Tokyo

# Course Plan

- Multi-core Resource Management
- Many-core Resource Management
- GPU Resource Management
- Virtual Machines
- Distributed File Systems
- High-performance Networking
- Memory Management
- Network on a Chip
- Embedded Real-time OS
- Device Drivers
- Linux Kernel

# Schedule

1. 2018.9.28 Introduction + Linux Kernel (Kato)
2. 2018.10.5 Linux Kernel (Chishiro)
3. 2018.10.12 Linux Kernel (Kato)
4. 2018.10.19 Linux Kernel (Kato)
5. 2018.10.26 Linux Kernel (Kato)
6. 2018.11.2 Advanced Research (Chishiro)
7. 2018.11.9 Advanced Research (Chishiro)
8. 2018.11.16 (No Class)
9. 2018.11.23 (Holiday)
10. 2018.11.30 Advanced Research (Kato)
11. 2018.12.7 Advanced Research (Kato)
12. 2018.12.14 Advanced Research (Chishiro)
13. 2018.12.21 Linux Kernel
14. 2019.1.11 Linux Kernel
15. 2019.1.18 (No Class)
16. 2019.1.25 (No Class)

# Process Management

Abstracting Computing Resources
/* The case for Linux */

*Acknowledgement:*

# Outline1

# Outline1

# Process

Definition

- e Refers to a **program currently executing in the system**
  - › CPU registers
  - › Location and state of memory segments (text, data, stack, etc.)
  - › Kernel resources (open files, pending signals, etc.)
  - › Threads
- e Managed on a per-program way:
  - › *Virtualization* of the processor and the memory
- e Let's check an example with `strace` (-f)

```
fork()
 clone()
                    exec()
                     execve()

                    exit()
                     exit()
wait()
 wait4()

          exit()
           exit_group()

parent   child
```

# Process

Sample program

```c
/* process.c */

#include <stdio.h>
#include <stdlib.h>
#include <unistd.h>
#include <sys/types.h>
#include <sys/wait.h>

int main(void)
{
  pid_t pid = -42;
  int wstatus = -42;
  int ret = -1;

  pid = fork();
  switch(pid)
  {
    case -1:
      perror("fork");
      return EXIT_FAILURE;

    case 0:
      sleep(1);
      printf("Noooooooo!¥n");
      exit(0);
```

```c
    default:
      printf("Iamyourfather!¥n");
      break;
  }

  ret = waitpid(pid, &wstatus, 0);
  if(ret == -1)
  {
    perror("waitpid");
    return EXIT_FAILURE;
  }
  printf("Childexitstatus:%d¥n",
      WEXITSTATUS(wstatus));

  return EXIT_SUCCESS;
}
```

```
gcc -Wall -Werror process.c -o process
./process
strace -f ./process > /dev/null
```

# Process

`fork()` & `exec()` usage

e Tutorial on `fork()` usage:

› http://www.csl.mtu.edu/cs4411.ck/www/NOTES/
process/fork/create.html

e Combining `fork()` and `exec()`:

› https://ece.uwaterloo.ca/dwharder/icsrts/
Tutorials/fork_exec/

# Outline1

# The process descriptor: `task_struct`

Presentation

- List of processes implemented as a linked list of `task_struct`

```
1  struct tastk_struct {
2    volatilelong state;
3    void *stack;
4    /* ... */
5    int prio;
6    /* ... */
7    cpumask_t cpus_allowed;
8    /* ... */
9    struct list_head tasks;
10   /* ... */
11   struct mm_struct *mm;
12   /* ... */
13   pid_t pid;
14   /* ... */
15   struct task_struct *parent;
16   struct list_head children;
17   struct list_head sibling;
18   /* ... */
19 }
```

- Total size (Linux 4.8): 6976 bytes
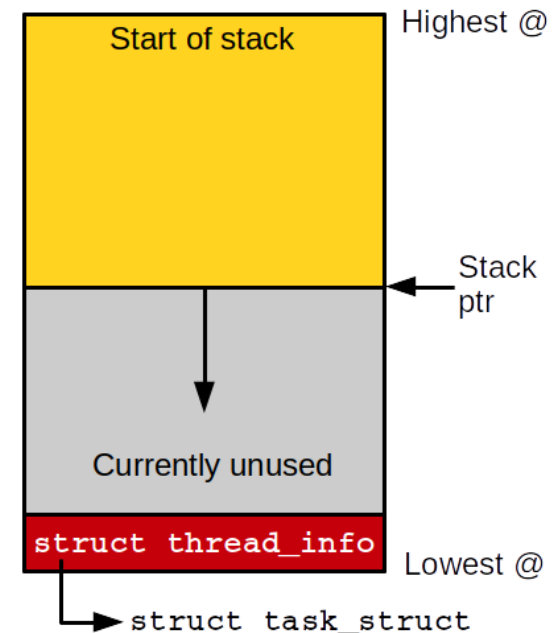- Full structure definition in `linux/sched.h`

# The process descriptor: `task_struct`

Allocation & storage

- e Prior to 2.6: `task_struct` allocated at the end of the kernel stack of each process
  - › Allows to retrieve it without storing its location in a register
- e Now dynamically allocated (heap) through the *slab allocator*
  - › A `struct thread_info` living at the bottom of the stack

```
1  struct thread_info {
2          struct task_struct        *task;
3          u32                       flags;
4          u32                       status;
5          u32                       cpu;
6  };
```



Start of stack — Highest @

Stack ptr

Currently unused

struct thread_info — Lowest @

struct task_struct

- e Moved off the stack in 4.9 [2] because of potential exploit [1] when overflowing the kernel stack

# The process descriptor: `task_struct`

Allocation & storage (2)

- e **Process Identifier (PID):** `pid_t` (`int`)
  - › Max: 32768, can be increased to 4 millions
  - › Wraps around when maximum reached
- e Quick access to `task_struct` of the task currently running on a core: `current`
  - › `arch/x86/include/asm/current.h`:

```
1
2  DECLARE_PER_CPU(struct task_struct *, current_task);
3
4  static_always_inline struct task_struct *get_current(void)
5  {
6      return this_cpu_read_stable(current_task);
7  }
8
9  #define current get_current()
```

# The process descriptor: `task_struct`

Process states

- e **state** field of the `task_struct`
  - ) **TASK_RUNNING:**
    - ) Process is runnable (running or in a CPU run queue)
    - ) In user or kernel space
  - ) **TASK_INTERRUPTIBLE:**
    - ) Process is sleeping waiting for some condition
    - ) Switched to `TASK_RUNNING` on condition true or signal received
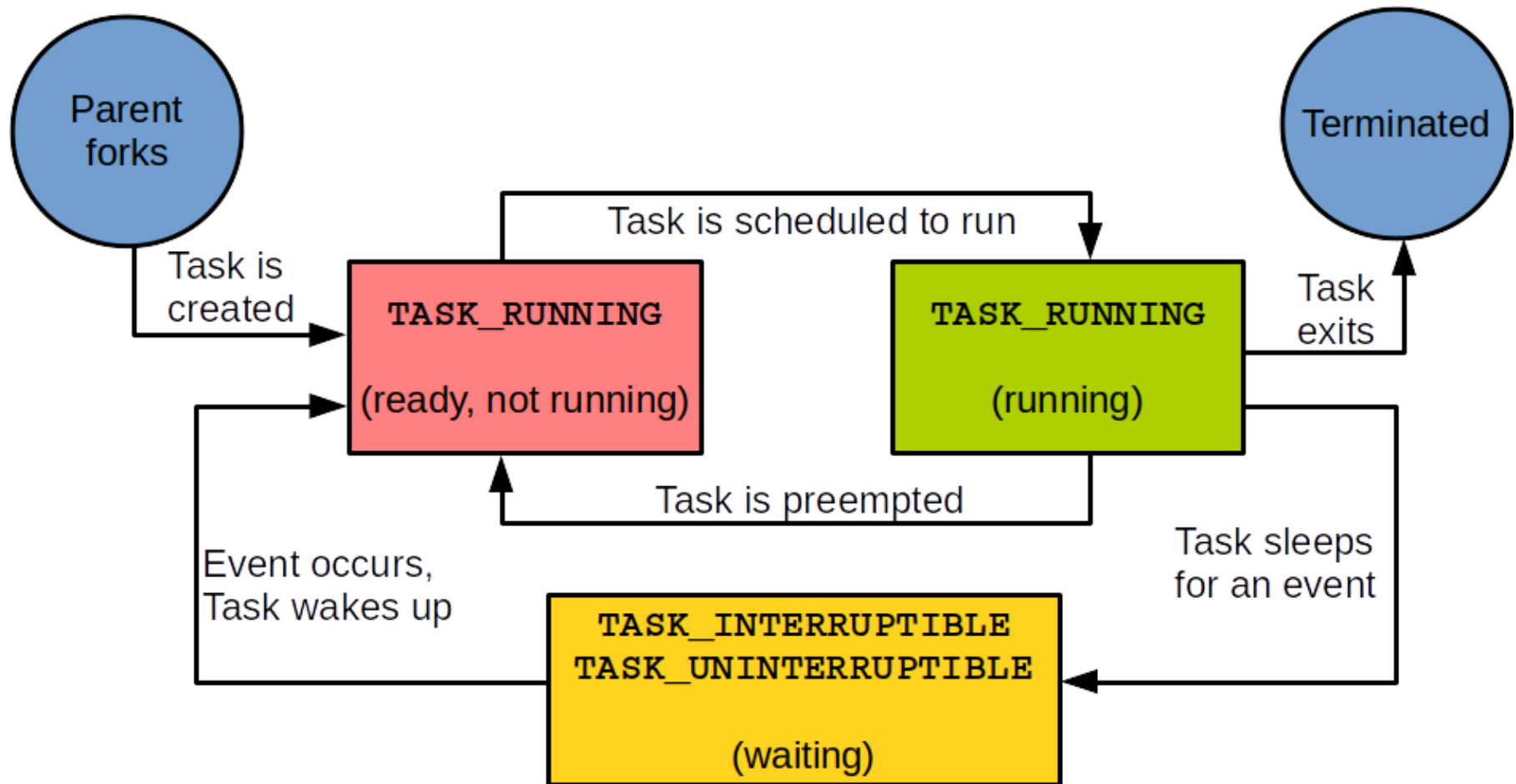  - ) **TASK_UNINTERRUPTIBLE:**
    - ) Same as `TASK_INTERRUPTIBLE` but does not wake up on signal
  - ) `TASK_TRACED`: Traced by another process (ex: debugger)
  - ) `TASK_STOPPED`: Not running nor waiting, result of the reception of some signals to pause the process

# The process descriptor: `task_struct`

Process states: flowchart

# The process descriptor: `task_struct`

Process context and family tree

- e The kernel can executes in **process** vs **interrupt** context
  - › `current` is meaningful only when the kernel executes in *process context*
    - › I.e. following a system call or an exception
- e **Process hierarchy**
  - › Root: `init`, PID 1
    - › Launched by the kernel as the last step of the boot process
  - › `fork`-based process creation:
    - › Each process has a parent: `parent` pointer in the `task_struct`
    - › Processes may have children: `children` field (`list_head`)
    - › Processes may have siblings: `siblings` field
    - › List of all tasks: `tasks` field
      - Easy manipulation through `next_task(t)` and `for_each_process(t)`
  - › Let's check it out with the `pstree` command

# Outline1

# Process creation

Presentation, Copy-On-Write

- e Linux does not implements creating a tasks from nothing (*spawn*)
- e **`fork()` & `exec()`**
    - › `fork()` creates a child, copy of the parent process
        - › Only PID, PPID and some resources/stats differ
    - › `exec()` loads into a process address space a new executable
- e On `fork()`, Linux duplicates the parent page tables and creates a new process descriptor
    - › It's *fast*, as **the address space is not copied**
        - › Page table access bits: read-only
        - › **Copy-On-Write** (COW): memory pages are copied only when they are referenced for write operations

# Process creation
Forking: `fork()` and `vfork()`

e `fork()` is implemented by the **`clone()`** system call

1. `sys_clone()` calls `do_fork()`, which calls **`copy_process()`** and starts the new task
2. **`copy_process()`**:
    1. Calls `dup_task_struct()`
        - Duplicates kernel stack, `task_struct` and `thread_info`
    2. Checks that we do not overflow the processes number limit
    3. Small amount of values are modified in the `task_struct`
    4. Calls `sched_fork()` to set the child `state` set to `TASK_NEW`
    5. Copies parent info: files, signal handlers, etc.
    6. Gets a new PID through `alloc_pid()`
    7. Returns a pointer to the created child `task_struct`
3. Finally, `do_fork()` calls `wake_up_new_task()`
    - State becomes `TASK_RUNNING`

e `vfork()`: alternative without copy of the address space _

# Outline1

# Threads
Presentation

e  Theory:



Execution flow — Address space

Time

A

Cpu core →

1 process, 1 thread

1 process, 3 threads

B C D

e **Threads** are concurrent flows of execution belonging to the same program **sharing the same address space**

e In Linux there is no concept of a thread
  › No scheduling particularity
  › A thread is just another process sharing some information with other processes
  › Each thread has its own `task_struct`
  › Created through `clone()` with specific flags indicating sharing

# Threads

Kernel threads

- To perform background operations in the kernel: **kernel threads**
- Very similar to user space threads
  - › They are *schedulable entities* (like regular processes)
- However **they do not have their own address space**
  - › `mm` in `task_struct` is `NULL`
- Used for several tasks:
  - › Work queues (`kworker`)
  - › Load balancing between CPU scheduling runqueues (`migration`)
  - › etc.
  - › List of all them with `ps --ppid 2`

# Threads

Kernel threads: creation

- e Kernel threads are all forked from the `kthread` kernel thread (PID 2), using `clone()`
    - › To create a kernel thread, use `kthread_create()`
    - › `include/linux/kthread.h`:

```
1  #define kthread_create(threadfn, data, namefmt, arg...) ¥
2    kthread_create_on_node(threadfn, data, NUMA_NO_NODE, namefmt, ##arg)
```

```
1  struct task_struct *kthread_create_on_node(int (*threadfn)(void *data),
2            void *data,
3            int node,
4            constchar namefmt[], ...);
```

- e When created through `kthread_create()`, the thread is not in a runnable state
    - › Need to call `wake_up_process()`:

```
1  int wake_up_process(struct task_struct *p);
```

    - › Or use `kthread_run()`

# Threads

Kernel threads: creation (2)

e `kthread_run():`

```
1 #define kthread_run(threadfn, data, namefmt, ...)              ¥
2 ({                                                             ¥
3    struct task_struct * k                                      ¥
4      = kthread_create(threadfn, data, namefmt, ##_VA_ARGS_);   ¥
5    if (!IS_ERR( k))                                            ¥
6      wake_up_process( k);                                      ¥
7    k;                                                          ¥
8 })
```

e Thread termination:

› Thread runs until it calls `do_exit():`

```
1 void do_exit(long error_code) noreturn;
```

› Or until another part of the kernel calls `kthread_stop():`

```
1 int kthread_stop(struct task_struct *k);
```

# Outline1

# Process termination

Termination steps: `do exit()`

- e Termination on invoking the `exit()` system call
  - ) Can be implicitly inserted by the compiler on `return` from `main`
  - ) `sys_exit()` calls `do_exit()`
- e `do_exit()` (`kernel/exit.c`):
  1. Calls `exit_signals()` which set the `PF_EXITING` flag in the `task_struct`
  2. Set the exit code in the `exit_code` field of the `task_struct`
     - ) To be retrieved by the parent
  3. Calls `exit_mm()` to release the `mm_struct` for the task
     - ) If it is not shared with any other process, it is destroyed
  4. Calls `exit_sem()`: process dequeued from potential semaphores queues
  5. Calls `exit_fs()` and `exit_files()` to update accounting information
     - ) Potential data structures that are not used anymore are freed

# Process termination

Termination steps: `do_exit()` (2)

- e `do_exit()` (continued):
    - 6 Calls `exit_notify()`
        - › Sends signals to parent
        - › Reparent potential children
        - › Set the `exit_state` of the `task_struct` to `EXIT_ZOMBIE`
    - 7 Calls `do_task_dead()`
        - › Sets the `state` to `TASK_DEAD`
        - › Calls `_schedule()` and never returns
- e At that point, what is left is the `task_struct`, `thread_info` and kernel stack
    - › To provide information to the parent
    - › Parent notifies the kernel when everything can be freed

# Process termination

`task_struct` cleanup

- e Separated from the process of exiting because of the need to pass exit information to the parent
  - › `task_struct` must survive a little bit before being deallocated
    - › Until the parent grab the exit information through `wait4()`
- e Cleanup implemented in `release_task()` called from the `wait4()` implementation
  - › Remove the task from the task list
  - › Release and free remaining resources

# Process termination

Parentless tasks

- <span style="color:orange">e</span> **A parent exits before its child**
    - › Child must be *reparented*
        - › To another process in the current thread group ...
        - › ... or `init` if that fails
- <span style="color:orange">e</span> `exit_notify()` calls `forget_original_parent()`, that calls `find_new_reaper()`
    - › Returns the `task_struct` of another task in the thread group if it exists, otherwise the one from `init`
    - › Then, all the children of the currently dying task are reparented to the reaper

# Bibliography I

1 Exploiting stack overflow in the linux kernel.
   https://jon.oberheide.org/blog/2010/11/29/exploiting-stack-overflows-in-the-linux-kernel/.
   Accessed: 2017-01-23.

2 Security things in linux v4.9.
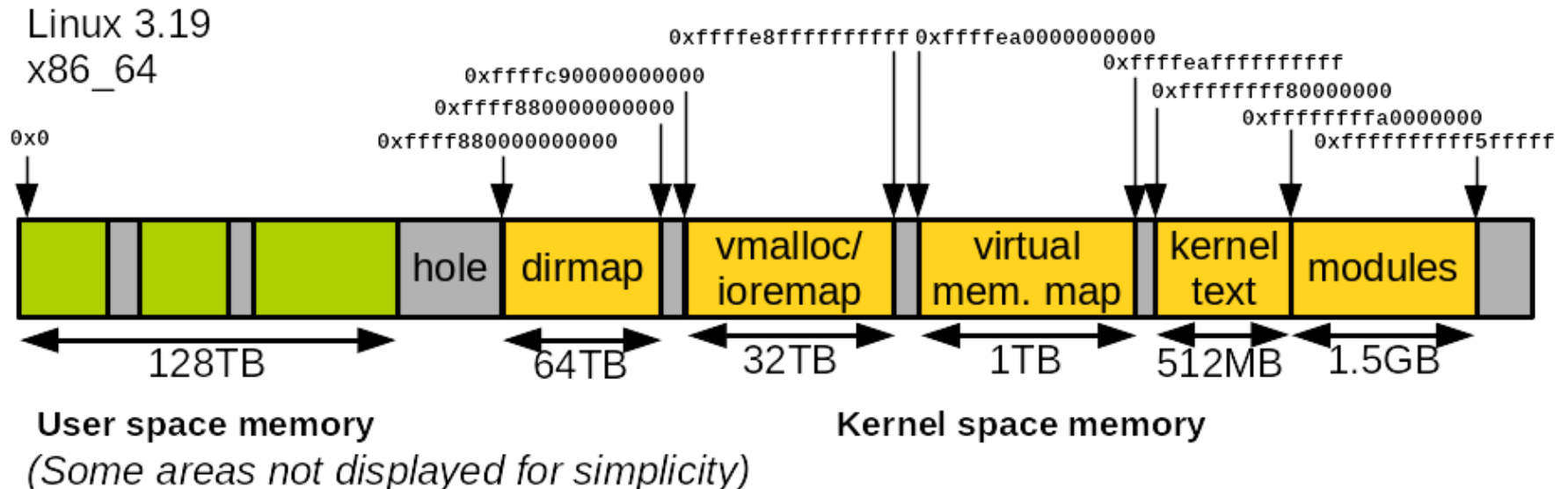   https://outflux.net/blog/archives/2016/12/12/security-things-in-linux-v4-9/.
   Accessed: 2017-01-23.

# Outline2

# Outline2

# Address space and memory descriptor

Address space

- e **The memory that a process can access is called its address space**
    - › Illusion that the process can access 100% of the system memory
    - › With virtual memory, can be much larger than the actual amount of physical memory
- e Defined by the process page table set up by the kernel



Linux 3.19 x86_64

0x0

0xffff880000000000
0xffff880000000000
0xffffc90000000000
0xffffe8ffffffffff 0xffffea0000000000
0xffffeaffffffffff
0xffffffff80000000
0xffffffffa0000000
0xffffffffff5fffff

| 128TB | hole | dirmap 64TB | vmalloc/ ioremap 32TB | virtual mem. map 1TB | kernel text 512MB | modules 1.5GB |

User space memory                    Kernel space memory
(Some areas not displayed for simplicity)

# Address space and memory descriptor

Address space (2)

e Each process is given a *flat* 32/64-bits address space
  › *Flat* as opposed to segmented



16-bit CPU

| Segment selector (base): `0x1000` |

| Address accessed: `0x1234` |

| Resulting physical @: `Base*F+offset` | → `0x10234` |

| Memory address (offset): `0x1234` |

Adapted from
http://duartes.org/gustavo/blog/post/memory-translation-and-segmentation/

e A **memory address** is an index within the address spaces:
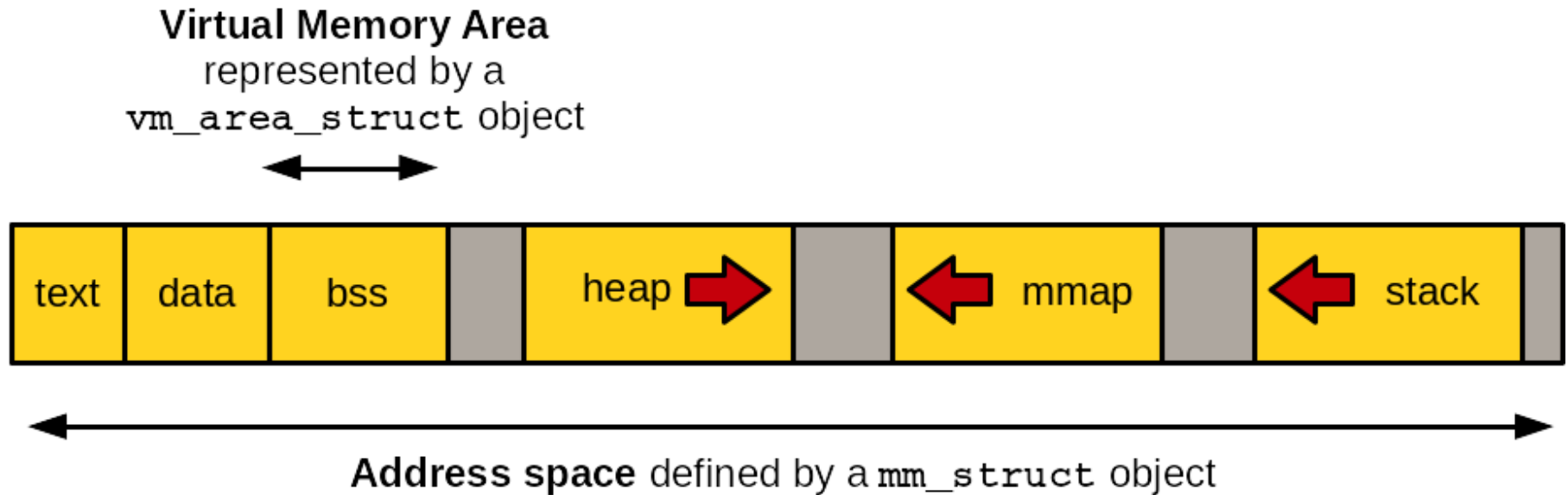  › Identify a specific byte
  › Example: `0x8fffa12dd24123fd`

# Address space and memory descriptor

Address space (3)

- Interval of addresses that the process has the right to access: **virtual memory areas** (**VMAs**)
  - VMAs can be dynamically added or removed to the process address space
  - VMAs have associated permissions: read, write, execute
    - When a process try to access an address outside of valid VMAs, or access a VMA with wrong permissions: **segmentation fault**
- **VMAs can contain**:
  - Mapping of the executable file code (*text section*)
  - Mapping of the executable file initialized variables (*data section*)
  - Mapping of the zero page for uninitialized variables (*bss section*)
  - Mapping of the zero page for the user-space stack
  - Text, data, bss for each shared library used
  - Memory-mapped files, shared memory segment, anonymous mappings (used by malloc)

# Address space and memory descriptor

Address space (4)

**Virtual Memory Area**
represented by a
`vm_area_struct` object

| text | data | bss | | heap | | mmap | | stack | |
|------|------|-----|--|------|--|------|--|-------|--|

**Address space** defined by a `mm_struct` object

# Address space and memory descriptor

Memory descriptor

e The kernel represent a process address space through a **`struct mm_struct`** object, the **memory descriptor**

› Defined in `include/linux/mm_types.h`
› Interesting fields:

```
 1  struct mm_struct {
 2   struct vm_area_struct *mmap;   /* listofVMAs        */
 3   struct rb_root        mm_rb;            /* rbtreeofVMAs     */
 4   pgd_t                 *pgd;             /* pageglobaldirectory */
 5   atomic_t              mm_users;         /* addressspaceusers */
 6   atomic_t              *mm_count;        /* primaryusagecounters */
 7   int                   map_count;        /* numberofVMAs       */
 8   struct rw_semaphore   mmap_sem;         /* VMAsemaphore */
 9   spinlock_t            page_table_lock;  /* pagetablelock    */
10   struct list_head      mmlist;           /* listofallmm_struct */
11   unsignedlong          start_code;       /* startaddressofcode */
12   unsignedlong          end_code;         /* endaddressofcode */
13   unsignedlong          start_data;       /* startaddressofdata */
14   unsignedlong          end_data;         /* endaddressofdata */
15  /* ... */
```

# Address space and memory descriptor

Memory descriptor (2)

```
16   /* ... */
17   unsignedlong            start_brk;  /* startaddressofheap    */
18   unsignedlong            end_brk;    /* endaddressofheap     */
19   unsignedlong            start_stack; /* startaddressofstack    */
20   unsignedlong            arg_start;  /* startofarguments    */
21   unsignedlong            arg_end;    /* endofarguments    */
22   unsignedlong            env_start;  /* startofenvironment    */
23   unsignedlong            total_vm;   /* totalpagesmapped    */
24   unsignedlong            locked_vm;  /* numberoflockedpages     */
25   unsignedlong            flags;      /* architecturespecificdata    */
26   spinlock_t              ioctx_lock; /* AsynchronousI/Olistlock     */
27   /* ... */
28 };
```

e `mm_users`: number of processes (threads) using the address space

e `mm_count`: reference count:
  ) +1 if `mm_users` > 0
  ) +1 if the kernel is using the address space
  ) When `mm_count` reaches 0, the `mm_struct` can be freed

# Address space and memory descriptor

Memory descriptor (3)

- e `mmap` and `mm_rb` are respectively a linked list and a tree containing all the VMAs in this address space
    - › List used to iterate over all the VMAs
        - › Links all VMAs sorted by ascending virtual addresses
    - › Tree used to find a specific VMA
- e All `mm_struct` are linked together in a doubly linked list
    - › Through the `mmlist` field if the `mm_struct`

# Address space and memory descriptor

Memory descriptor allocation

- A task memory descriptor is located in the `mm` field of the corresponding `task_struct`
  - Current task memory descriptor: `current->mm`
  - During `fork()`, `copy_mm()` is making a copy of the parent memory descriptor for the child
    - `copy_mm()` calls `dup_mm()` which calls `allocate_mm()` → allocates a `mm_struct` object from a slab cache
  - Two threads sharing the same address space have the `mm` field of their `task_struct` pointing to the same `mm_struct` object
    - Threads are created using the `CLONE_VM` flag passed to `clone()` → `allocate_mm()` is not called
    - in `copy_mm()`:

```
1  /* ... */
2  struct mm_struct *oldmm;
3  oldmm = current->mm;
4  /* ... */
5  if (clone_flags & CLONE_VM) {
6    atomic_inc(&oldmm->mm_users);
7    mm = oldmm;
```

```
8    goto good_mm;
9  }
10 /* ... */
11 good_mm:
12   /* ... */
13 return 0;
```

# Address space and memory descriptor

Memory descriptor destruction

- When a process exits, `do_exit()` is called
  - It calls `exit_mm()`
    - Performs some housekeeping/statistics updates
    - Calls `mmput()`

```
1  void mmput(struct mm_struct *mm)
2  {
3    might_sleep();
4
5    if (atomic_dec_and_test(&mm->mm_users))
6      __mmput(mm);
7  }
```

  - `mmput()` decrements the `users` field and calls`_mmput()` if it reaches 0
  - `_mmput()` calls `mmdrop()`, that decrements the `count` field, and calls`_mm_drop()` if it reaches 0
  - `_mmdrop()` calls `free_mm()` which return the memory for the `mm_struct()` to the slab cache (i.e. free)

# Address space and memory descriptor

Memory descriptor and kernel threads



e Kernel threads do not have a user-space address space

 › `mm` field of a kernel thread `task_struct` is `NULL`

e However they still need to access the kernel address space
 › When a kernel thread is scheduled, the kernel notice its `mm` is `NULL` so it keeps the previous address space loaded (page tables)
 › Kernel makes the `active_mm` field of the kernel thread to point on the *borrowed* `mm struct`
 › OK because **kernel part is the same in all address spaces**

# Outline2

# Virtual Memory Area

vm_area_srtuct

e Each VMA is represented by an object of type **vm area struct**

  › Defined in `include/linux/mm_types.h`
  › Interesting fields:

```
 1 struct vm_area_struct {
 2     struct i                      mm_struct *vm_mm; /* associatedaddressspace(mm_struct)    */
 3     unsignedlong                  vm_start;        /* VMAstart,inclusive   */
 4     unsignedlong                  vm_end;          /* VMAend,exclusive    */
 5     struct vm_area_struct         *vm_next;        /* listofVMAs    */
 6     struct vm_area_struct         *vm_prev;        /* listofVMAs    */
 7     pgprot_t                      vm_page_prot;    /* accesspermissions  */
 8     unsignedlong                  vm_flags;        /* flags */
 9     struct rb_node                vm_rb;           /* VMAnodeinthetree      */
10     struct list_head              anon_vma_chain;  /* listofanonymousmappings     */
11     struct anon_vma               *anon_vma;       /* anonmousvmaobject    */
12     struct vm_operation_struct *vm_ops;            /* operations */
13     unsignedlong                  vm_pgoff;        /* offsetwithinfile    */
14     struct file                   *vm_file;        /* mappedfile(canbeNULL)     */
15     void                          *vm_private_data; /* privatedata */
16     /* ... */
17 }
```

# Virtual Memory Area

`vm_area_srtuct` (2)

- e  The VMA exists over `[vm_start, vm_end[` in the corresponding address space
  - › Address space is pointed by the `vm_mm` field (of type `mm_struct`)
- e  Size in bytes: `vm_end - vm_start`
- e  Each VMA is unique to the associated `mm_struct`
  - › Two processes mapping the same file will have two different `mm_struct` objects, and two different `vm_area_struct` objects
  - › Two threads sharing a `mm_struct` object also share the `vm_area_struct` objects

# Virtual Memory Area

Flags

- **Flags** specify properties and information for all the pages contained in the VMA

- `VM_READ`: pages can be read from
- `VM_WRITE`: pages can be written to
- `VM_EXEC`: code inside pages can be executed
- `VM_SHARED`: pages are shared between multiple processes (if unset the mapping is *private*)
- `VM_MAYREAD`: the `VM_READ` flag can be set
- `VM_MAYWRITE`: the `VM_WRITE` flag can be set
- `VM_MAYEXEC`: the `VM_EXEC` flag can be set

- `VM_MAYSHARE`: the `VM_SHARED` flag can be set
- `VM_GROWSDOWN`: area can grow downwards
- `VM_GROWSUP`: area can grow upwards
- `VM_SHM`: area can be used for shared memory
- `VM_DENYWRITE`: area maps an unwritable file
- `VM_EXECUTABLE`: area maps an executable file

# Virtual Memory Area
Flags (2)

- `VM_LOCKED`: the area pages are locked (will not be swapped-out)
- `VM_IO`: the area maps a device IO space
- `VM_SEQ_READ`: pages in the area seem to be accessed sequentially
- `VM_RAND_READ`: pages seem to be accessed randomly

- `VM_DONTCOPY`: area will not be copied upon `fork()`
- `VM_DONTEXPAND`: area cannot grow through `mremap()`
- `VM_ACCOUNT`: area is an accounted VM object
- `VM_HUGETLB`: area uses `hugetlb` pages

# Virtual Memory Area

Flags (3)

- e Combining `VM_READ`, `VM_WRITE` and `VM_EXEC` gives the permissions for the entire area, for example:
  - › Object code is `VM_READ` and `VM_EXEC`
  - › Stack is `VM_READ` and `VM_WRITE`
- e `VM_SEQ_READ` and `VM_RAND_READ` are set through the `madvise()` system call
  - › Instructs the file pre-fetching algorithm *read-ahead* to increase or decrease its agressivity
- e `VM_HUGETLB` indicates that the area uses pages larger than the regular size
  - › 2M and 1G on x86
  - › Smaller page table → good for the TLB
  - › Less levels of page tables → faster address translation

# Virtual Memory Area

VMA operations

- e `vm_ops` in `vm_area_struct` points to a **vm_operations_struct** object
  - › Contains function pointers to operate on a specific VMAs
  - › Defined in `include/linux/mm.h`

```c
struct vm_operations_struct {
  void (*open)(struct vm_area_struct * area);
  void (*close)(struct vm_area_struct * area);
  int (*fault)(struct vm_area_struct *vma, struct vm_fault *vmf);
  int (*page_mkwrite)(struct vm_area_struct *vma, struct vm_fault *vmf);
  /* ... */
}
```
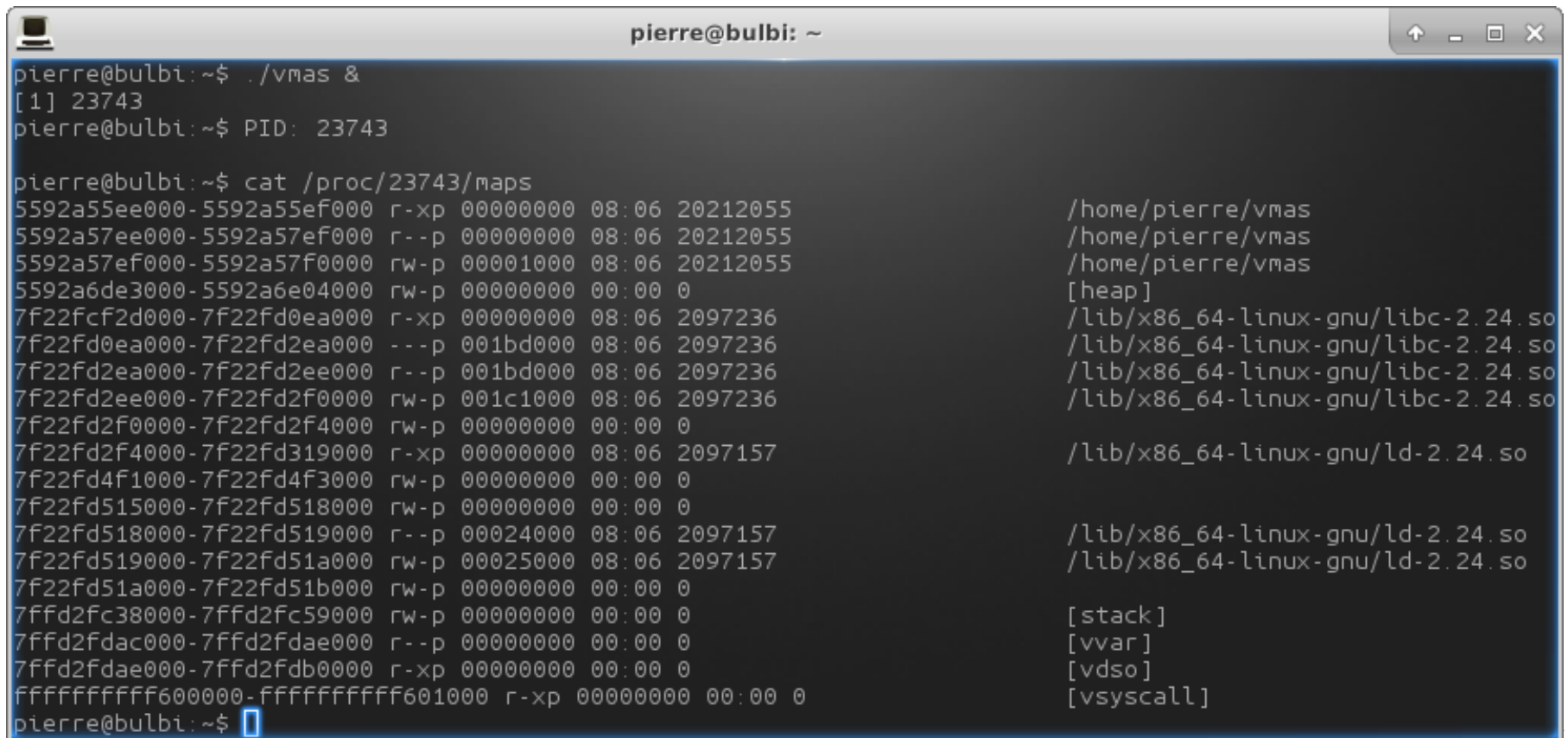
# Virtual Memory Area

VMA operations (2)

e Function pointers in `vm_operations_struct`:

> `open()`: called when the area is added to an address space
> `close()`: called when the area is removed from an address space
> `fault()`: invoked by the page fault handler when a page that is not present in physical memory is accessed
> `page_mkwrite()`: invoked by the page fault handler when a previously read-only page is made writable
> Description of all operations in `include/linux/mm.h`

# Virtual Memory Area

VMAs in real life

- e From userspace, one can observe the VMAs map for a given process:
  - › `cat /proc/<pid>/maps`

```
pierre@bulbi:~$ ./vmas &
[1] 23743
pierre@bulbi:~$ PID: 23743

pierre@bulbi:~$ cat /proc/23743/maps
5592a55ee000-5592a55ef000 r-xp 00000000 08:06 20212055         /home/pierre/vmas
5592a57ee000-5592a57ef000 r--p 00000000 08:06 20212055         /home/pierre/vmas
5592a57ef000-5592a57f0000 rw-p 00001000 08:06 20212055         /home/pierre/vmas
5592a6de3000-5592a6e04000 rw-p 00000000 00:00 0                [heap]
7f22fcf2d000-7f22fd0ea000 r-xp 00000000 08:06 2097236          /lib/x86_64-linux-gnu/libc-2.24.so
7f22fd0ea000-7f22fd2ea000 ---p 001bd000 08:06 2097236          /lib/x86_64-linux-gnu/libc-2.24.so
7f22fd2ea000-7f22fd2ee000 r--p 001bd000 08:06 2097236          /lib/x86_64-linux-gnu/libc-2.24.so
7f22fd2ee000-7f22fd2f0000 rw-p 001c1000 08:06 2097236          /lib/x86_64-linux-gnu/libc-2.24.so
7f22fd2f0000-7f22fd2f4000 rw-p 00000000 00:00 0
7f22fd2f4000-7f22fd319000 r-xp 00000000 08:06 2097157          /lib/x86_64-linux-gnu/ld-2.24.so
7f22fd4f1000-7f22fd4f3000 rw-p 00000000 00:00 0
7f22fd515000-7f22fd518000 rw-p 00000000 00:00 0
7f22fd518000-7f22fd519000 r--p 00024000 08:06 2097157          /lib/x86_64-linux-gnu/ld-2.24.so
7f22fd519000-7f22fd51a000 rw-p 00025000 08:06 2097157          /lib/x86_64-linux-gnu/ld-2.24.so
7f22fd51a000-7f22fd51b000 rw-p 00000000 00:00 0
7ffd2fc38000-7ffd2fc59000 rw-p 00000000 00:00 0                [stack]
7ffd2fdac000-7ffd2fdae000 r--p 00000000 00:00 0                [vvar]
7ffd2fdae000-7ffd2fdb0000 r-xp 00000000 00:00 0                [vdso]
ffffffffff600000-ffffffffff601000 r-xp 00000000 00:00 0        [vsyscall]
pierre@bulbi:~$
```

# Virtual Memory Area

VMAs in real life (2)

e `/proc/<pid>/maps` columns description:

1. Address range
2. Permissions
3. Start offset of file mapping
4. Device containing the mapped file
5. Mapped file inode number
6. Mapped file pathname

e Can also use the command `pmap <pid>`

# Outline2

# VMA manipulation

Finding a VMA

e **`find vma()`**: used to find a VMA in which a specific memory address resides

> Prototype in `include/linux/mm.h`:

```
1  struct vm_area_struct *find_vma(struct mm_struct *mm, unsignedlong addr);
```

> Defined in `mm/mmap.c`:

```
1  struct vm_area_struct *find_vma(struct
       mm_struct *mm, unsignedlong addr)
2  {
3    struct rb_node *rb_node;
4    struct vm_area_struct *vma;
5
6    /* Checkthecachefirst.    */
7    vma = vmacache_find(mm, addr);
8    if (likely(vma))
9      return vma;
10
11   rb_node = mm->mm_rb.rb_node;
12
13   while (rb_node) {
14     struct vm_area_struct *tmp;
```

```
15     tmp = rb_entry(rb_node, struct
         vm_area_struct, vm_rb);
16
17     if (tmp->vm_end > addr) {
18       vma = tmp;
19       if (tmp->vm_start <= addr)
20         break;
21       rb_node = rb_node->rb_left;
22     } else
23       rb_node = rb_node->rb_right;
24   }
25
26   if (vma)
27     vmacache_update(addr, vma);
28   return vma;
29 }
```

# VMA manipulation

Finding a VMA (2)

- e `find_vma_prev()`: returns in addition the last VMA before a given address
  - › `include/linux/mm.h`:

```
1  struct vm_area_struct *find_vma_prev(struct mm_struct *mm, unsigned long addr,
        struct vm_area_struct **pprev);
```

- e `find_vma_intersection()`: returns the first VMA overlapping a given address range
  - › `include/linux/mm.h`:

```
1  static inline struct vm_area_struct * find_vma_intersection(struct mm_struct *
        mm, unsigned long start_addr, unsigned long end_addr)
2  {
3    struct vm_area_struct * vma = find_vma(mm,start_addr);
4
5    if (vma && end_addr <= vma->vm_start)
6      vma = NULL;
7    return vma;
8  }
```

# VMA manipulation

Creating an address interval

- **`do_mmap()`** is used to create a new *linear address interval*:
  - › Can result in the creation of a new VMAs
  - › Or a merge of the create area with an adjacent one when they have the same permissions
  - › `include/linux/mm.h`:

```
1  externunsignedlong    do_mmap(struct file *file, unsignedlong addr,
2      unsignedlong len, unsignedlong prot, unsignedlong flags,
3      vm_flags_t vm_flags, unsignedlong pgoff, unsignedlong *populate);
```

- Caller must hold `mm->mmap_sem` (RW semaphore)

- Maps the file `file` in the address space at address `addr` for length `len`. Mapping starts at offset `pgoff` in the file

- `prot` specifies access permissions for the memory pages:
  - › `PROT_READ`, `PROT_WRITE`, `PROT_EXEC`, `PROT_NONE`

# VMA manipulation

Creating an address interval (2)

- e `flags` specifies the rest of the VMAoptions:

- e `MAP_SHARED`: mapping can be shared

- e `MAP_PRIVATE`: mapping cannot be shared

- e `MAP_FIXED`: created interval *must* start at `addr`

- e `MAP_ANONYMOUS`: mapping is not file-backed

- e `MAP_GROWSDOWN`: corresponds to `VM_GROWSDOWN`

- e `MAP_DENYWRITE`: corresponds to `VM_DENYWRITE`

- e `MAP_EXECUTABLE`: corresponds to `VM_EXECUTABLE`

- e `MAP_LOCKED`: corresponds to `VM_LOCKED`

- e `MAP_NORESERVE`: no space reserved for the mapping

- e `MAP_POPULATE`: populate (default) page tables

- e `MAP_NONBLOCK`: do not block on IO

# VMA manipulation

Creating an address interval (3)

- e On error `do_mmap()` returns a negative value
- e On success:
  - › The kernel tries to merge the new interval with an adjacent one having same permissions
  - › Otherwise, create a new VMA
  - › Returns a pointer to the start of the mapped memory area
- e `do_mmap()` is exported to user-space through `mmap2()`

```
1  void *mmap2(void *addr, size_t length, int prot, int flags, int fd, off_t pgoffset);
```

# VMA manipulation

Removing an address interval

- e Removing an address interval is done through `do munmap()`
    - › `include/linux/mm.h`:

    ```
    1 int do_munmap(struct mm_struct *, unsignedlong , size_t);
    ```

    - › 0 returned on success
- e Exported to user-space through `munmap()`:

```
1 int munmap(void *addr, size_t len);
```
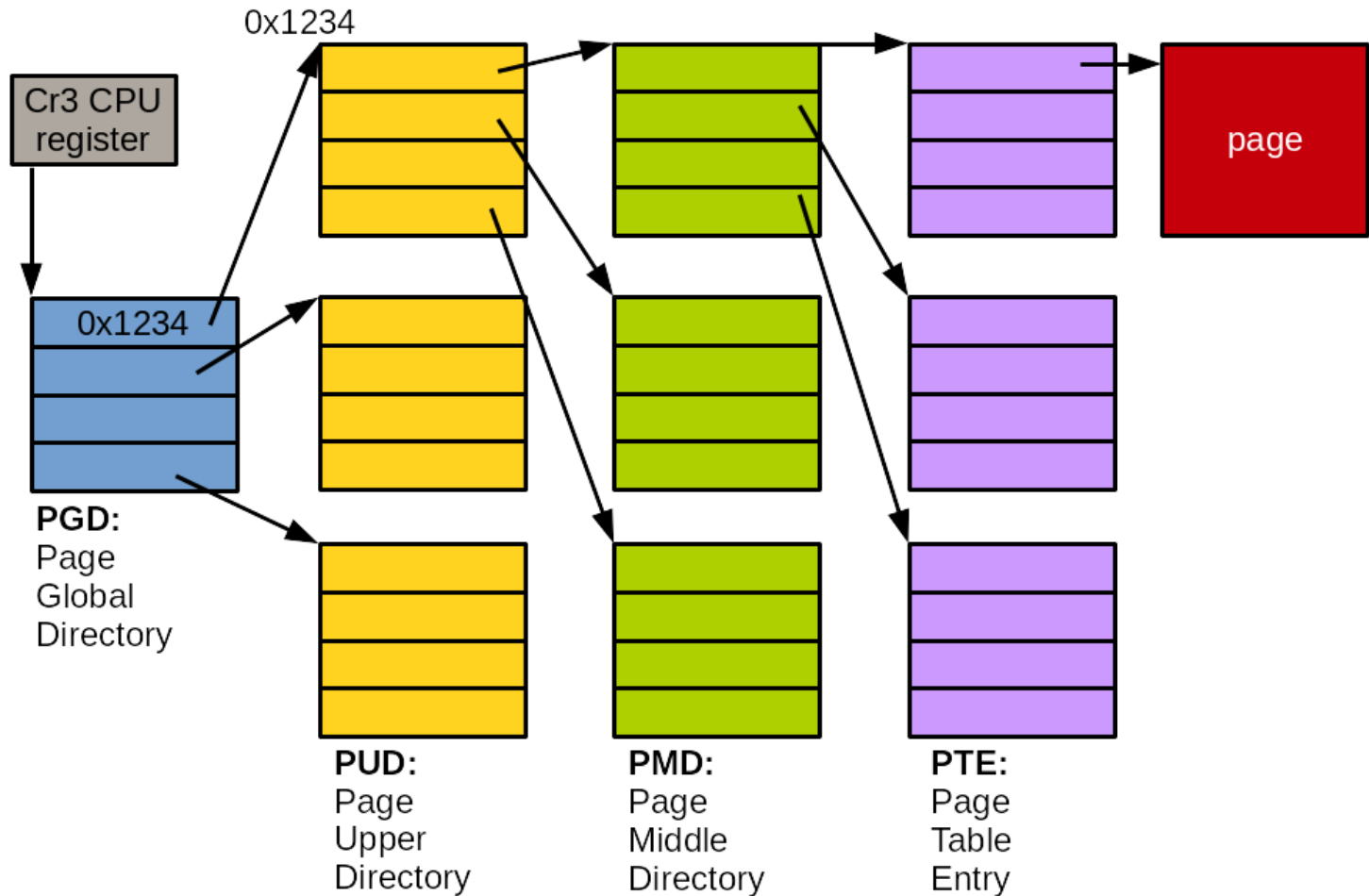
# Outline2

# Page tables

Presentation

e Linux enables **paging** early in the boot process

› **All memory accesses made by the CPU are virtual and translated to physical addresses through the *page tables***

› Linux set the page tables and the translation is made automatically by the hardware (MMU) according to the page tables content

e The address space is defined by VMAs and is sparsely populated

› One address space per process → one page table per process

› Lots of "empty" areas

› **Defining the page table as a single static array is a huge waste of space**

› A hierarchical tree structure is used
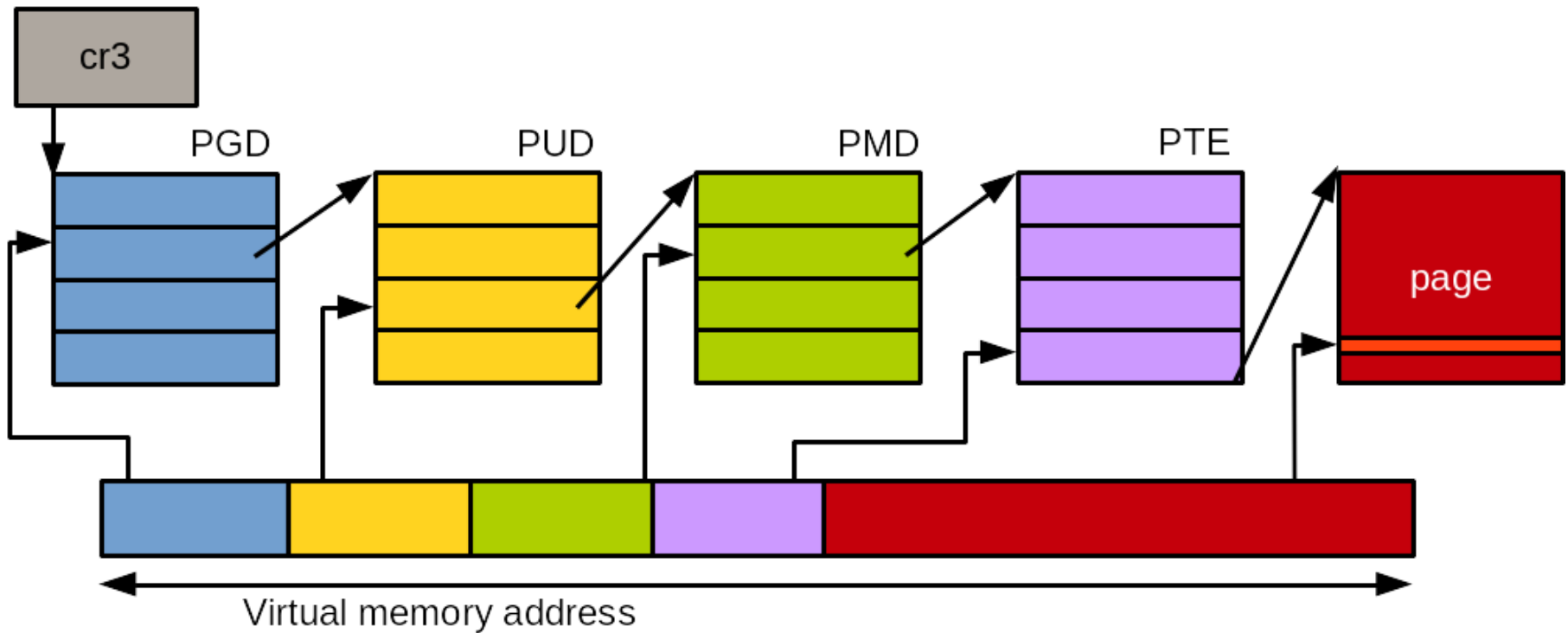
# Page tables

Page table setup (2)

e **Setting up the page table is performed by the kernel**

# Page tables

Address translation

- **Address translation is performed by the hardware**



- More info on page tables and memory management: [2, 1]

# BibliographyII

[1]Four-level page tables.
https://lwn.net/Articles/106177/.
Accessed: 2017-03-25.

[2]G ORMAN, M.
*Understanding the Linux virtual memory manager*.
Prentice Hall Upper Saddle River, 2004.
Accessed: 2017-03-25.