

# CS353 Linux Kernel

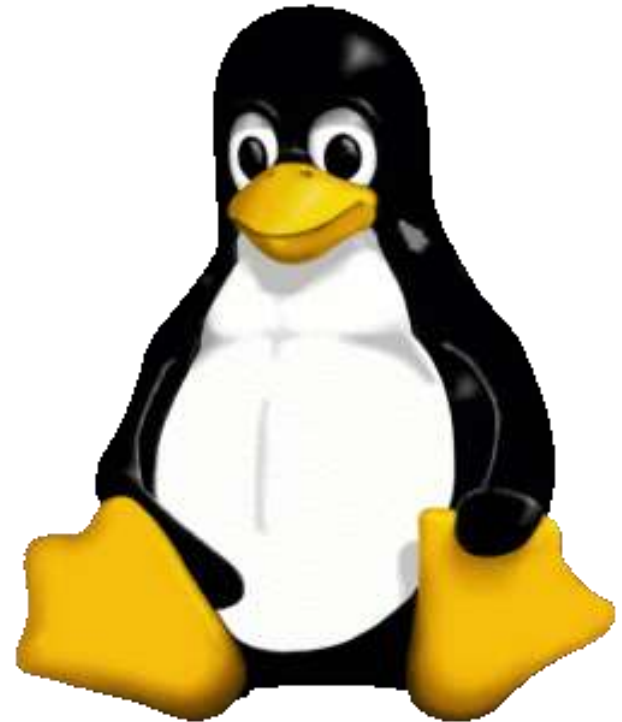
**Chentao Wu 吴晨涛**  
**Associate Professor**  
**Dept. of CSE, SJTU**  
**wuct@cs.sjtu.edu.cn**



上海交通大学

# 3A. Process Management -- Introduction

**Chentao Wu**  
Associate Professor  
Dept. of CSE, SJTU  
[wuct@cs.sjtu.edu.cn](mailto:wuct@cs.sjtu.edu.cn)



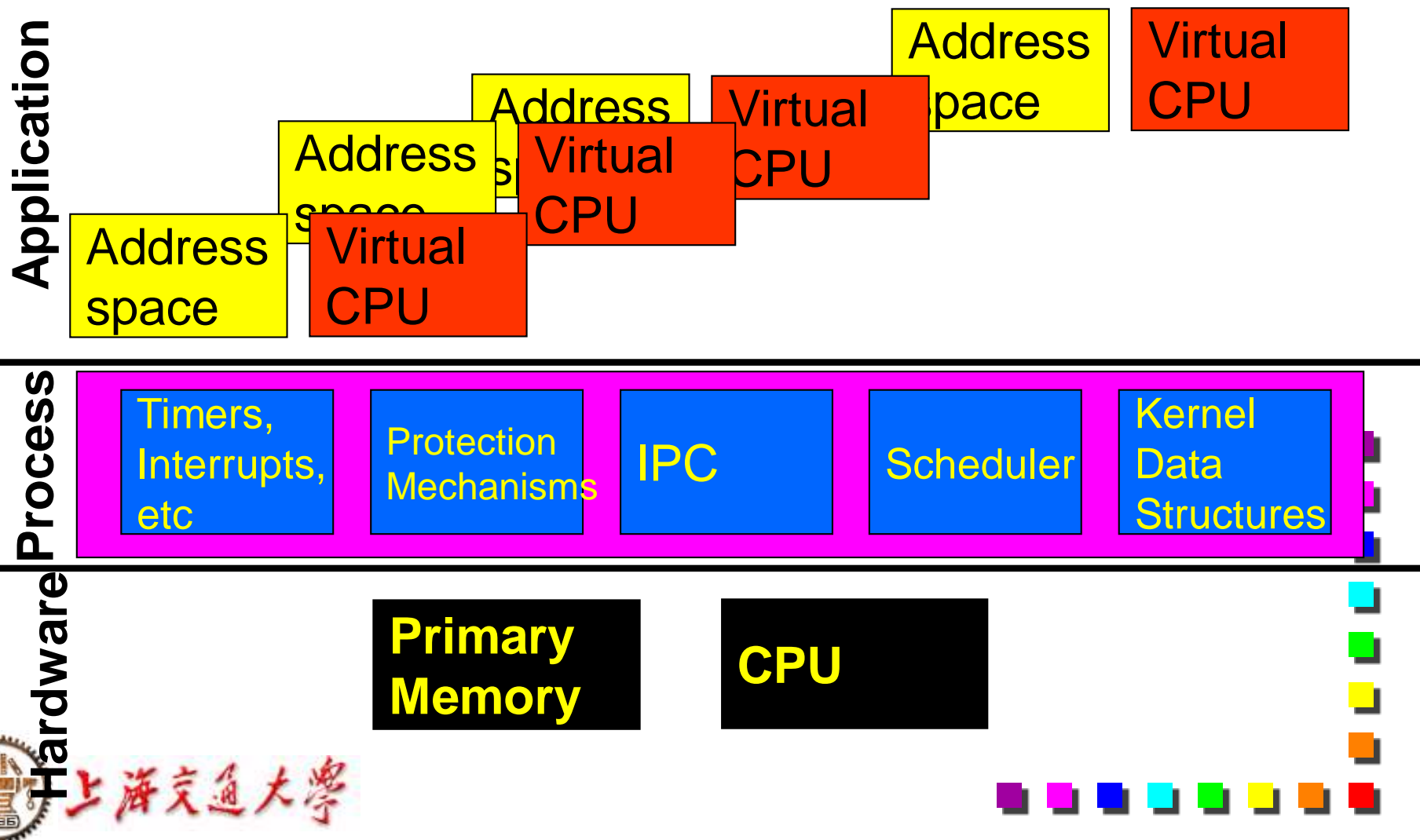
上海交通大學

# Processes, Lightweight Processes, and Threads

- *Process*: an instance of a program in execution
- *(User) Thread*: an execution flow of the process
  - Pthread (POSIX thread) library
- *Lightweight process (LWP)*: used to offer better support for multithreaded applications
  - LWP may share resources: address space, open files, ...
  - To associate a lightweight process with each thread
  - Examples of pthread libraries that use LWP: LinuxThreads, IBM's Next Generation Posix Threading Package (NGPT)

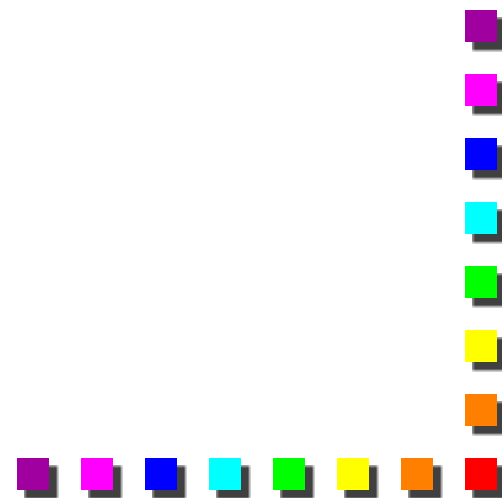


# Process Management

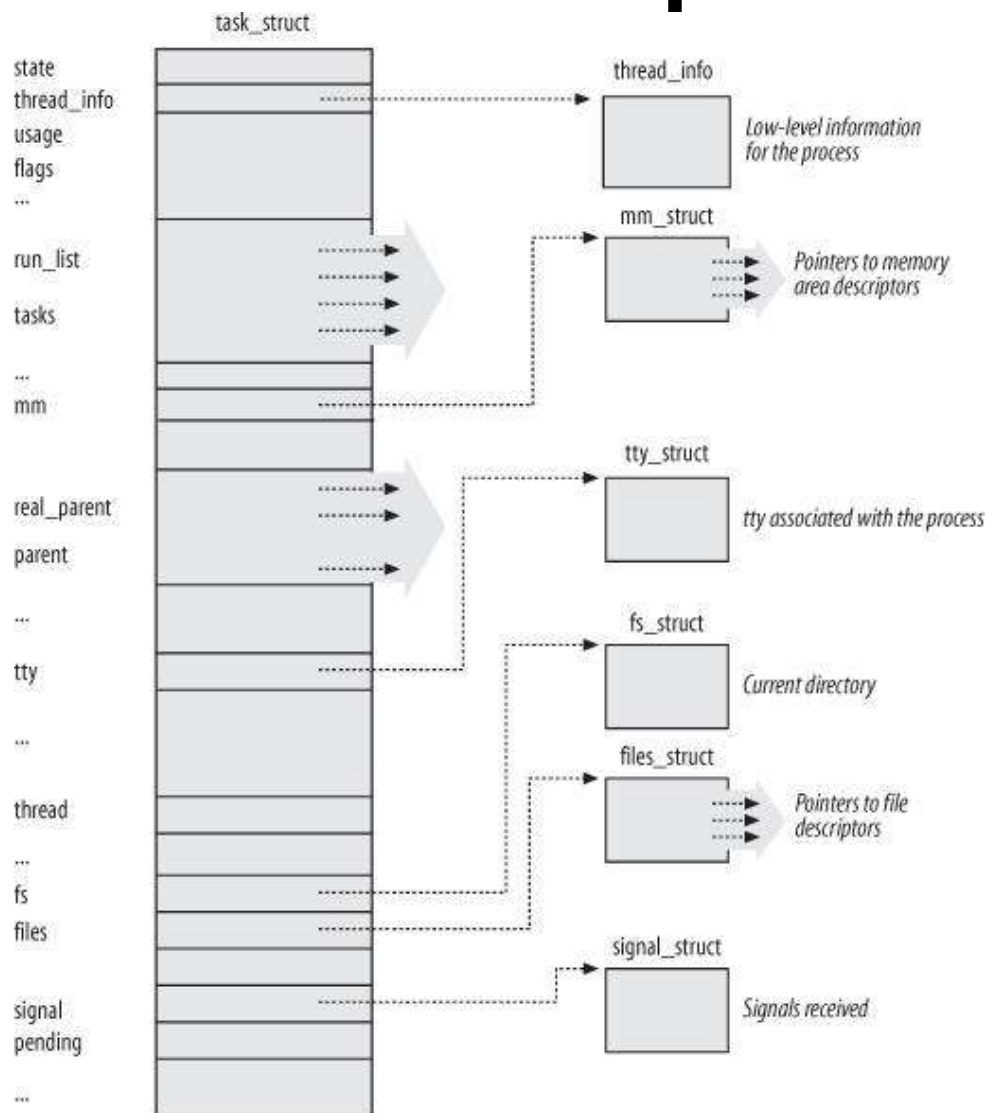


# Process Descriptor

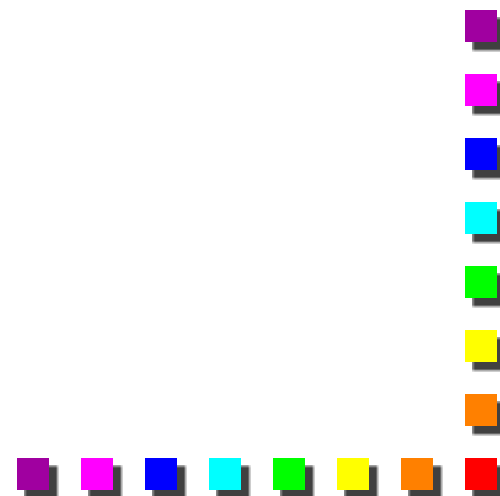
- *task\_struct* data structure
  - state: process state
  - thread\_info: low-level information for the process
  - mm: pointers to memory area descriptors
  - tty: tty associated with the process
  - fs: current directory
  - files: pointers to file descriptors
  - signal: signals received
  - ...



# Linux Process Descriptor



上海交通大学



# Process State

- TASK\_RUNNING: executing
- TASK\_INTERRUPTIBLE: suspended (sleeping)
- TASK\_UNINTERRUPTIBLE: (seldom used)
- TASK\_STOPPED
- TASK\_TRACED
- EXIT\_ZOMBIE
- EXIT\_DEAD



# Identifying a Process

- Process descriptor pointers: 32-bit
- Process ID (PID): 16-bit (~32767 for compatibility)
  - Linux associates different PID with each process or LWP
  - Programmers expect threads in the same group to have a common PID
  - *Thread group*: a collection of LWPs (kernel 2.4)
    - The PID of the first LWP in the group
    - *tgid* field in process descriptor: using getpid() system call



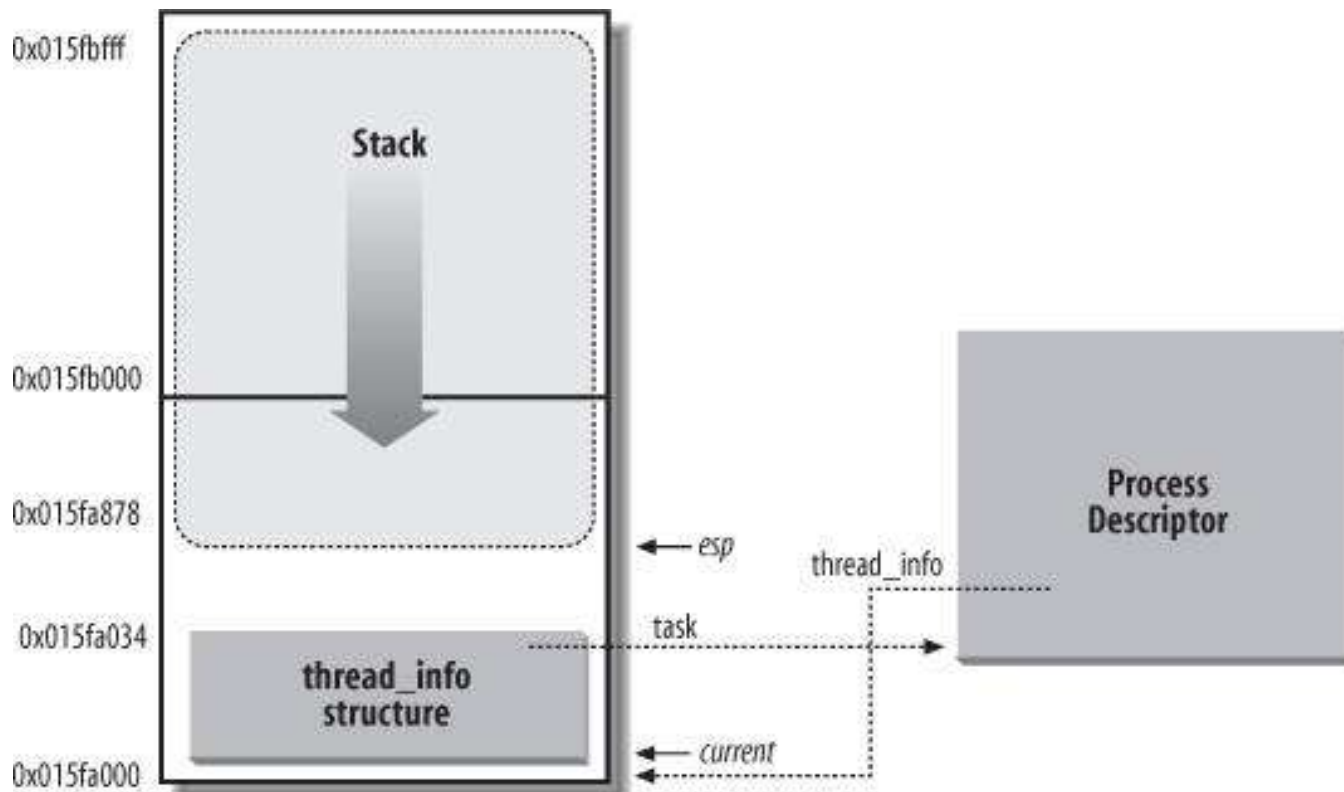


# Process Descriptor Handling

- union thread\_union {  
    struct thread\_info thread\_info;  
    unsigned long stack[2048];  
};
- Two data structures in 8KB (2 pages)
  - *thread\_info* structure (new to 3rd ed.)
  - Kernel mode process stack

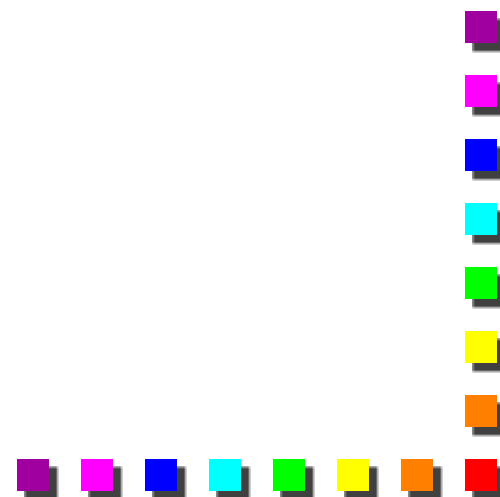


# Storing the *thread\_info* Structure and the Process Kernel Stack

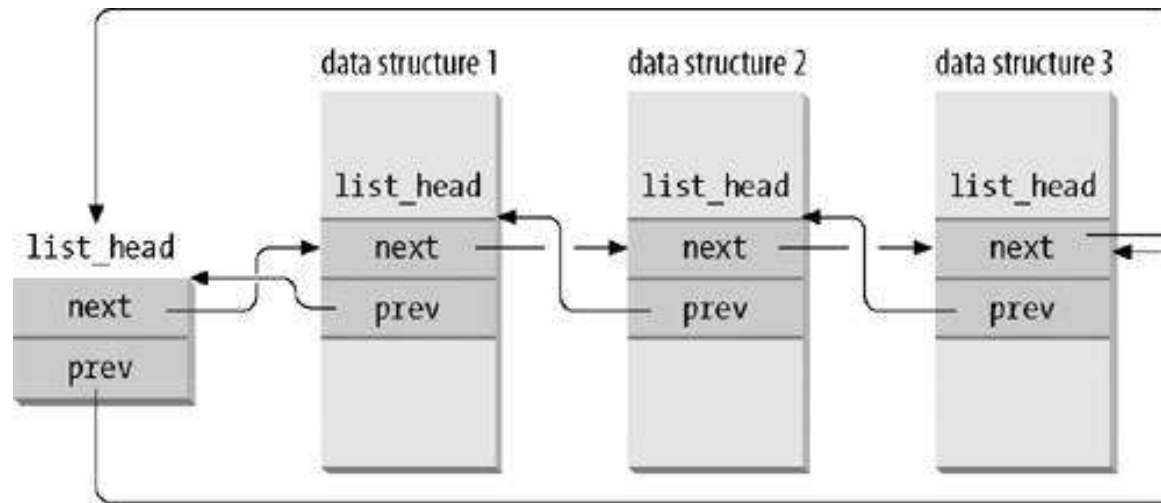


# Identifying the Current Process

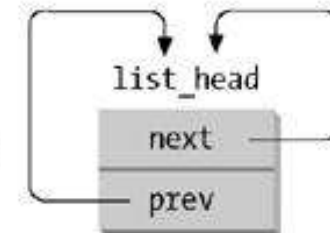
- Obtain the address of *thread\_info* structure from the *esp* register
  - `current_thread_info()`
  - `movl $0xffffe000, %ecx`  
`andl %esp, %ecx`  
`movl %ecx, p`
- The process descriptor pointer of the process currently running on a CPU
  - The *current* macro: equivalent to `current_thread_info()->task`
  - `movl $0xffffe000, %ecx`  
`andl %esp, %ecx`  
`movl (%ecx), p`



# Doubly Linked Lists Built with the *list\_head* Data Structure



(a) a doubly linked list with three elements



(b) an empty doubly linked list



# List Handling Functions and Macros

- LIST\_HEAD(list\_name)
- list\_add(n,p)
- list\_add\_tail(n,p)
- list\_del(p)
- list\_empty(p)
- list\_entry(p,t,m)
- list\_for\_each(p,h)
- list\_for\_each\_entry(p,h,m)



# The Process List

- *tasks* field in *task\_struct* structure
  - type *list\_head*
  - *prev*, *next* fields point to the previous and the next *task\_struct*
- Process 0 (swapper): *init\_task*
- Useful macros:
  - SET\_LINKS, REMOVE\_LINKS: insert and remove a process descriptor
  - #define for\_each\_process(p) \for (p=&init\_task; (p=list\_entry((p)->tasks.next), \struct task\_struct, tasks) \) != &init\_task; )



# List of TASK\_RUNNING processes

- *runqueue*
  - *run\_list* field in task\_struct structure: type list\_head
- Linux 2.6 implements the runqueue differently
  - To achieve scheduler speedup, Linux 2.6 splits the runqueue into 140 lists of each priority!
  - *array* field of process descriptor: pointer to the *prio\_array\_t* data structure
    - nr\_active: # of process descriptors in the list
    - bitmap: priority bitmap
    - queue: the 140 list\_heads
  - enqueue\_task(p, array), dequeue\_task(p, array)



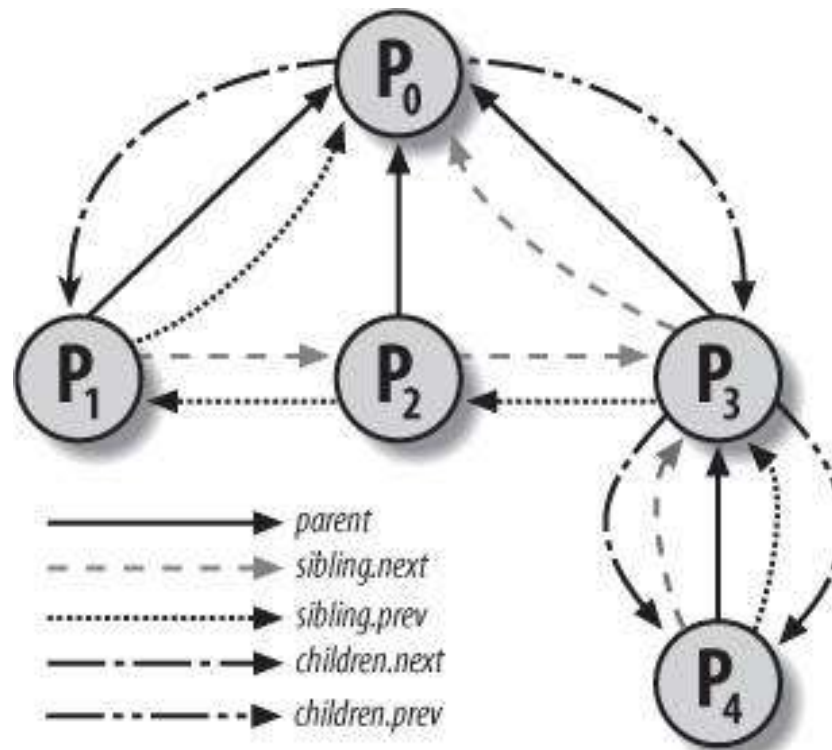
# Parenthood Relationships among Processes

- Process 0 and 1: created by the kernel
  - Process 1 (*init*): the ancestor of all processes
- Fields in process descriptor for parenthood relationships
  - `real_parent`
  - `parent`
  - `children`
  - `sibling`





# Parenthood Relationships among Five Processes



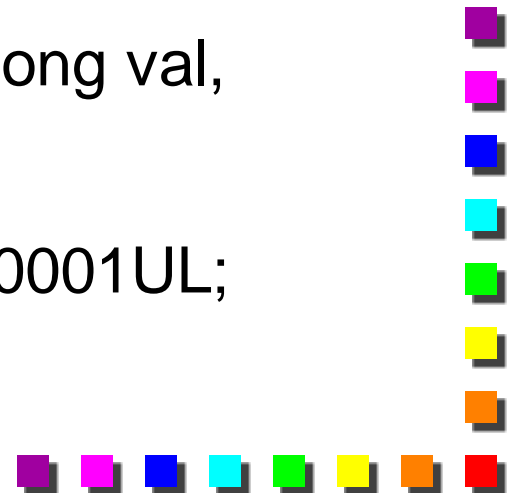
# Pidhash Table and Chained Lists

- To search up the search for the process descriptor of a PID
  - Sequential search in the process list is inefficient
- The *pid\_hash* array contains four hash tables and corresponding filed in the process descriptor
  - pid: PIDTYPE\_PID
  - tgid: PIDTYPE\_TGID (thread group leader)
  - pgrp: PIDTYPE\_PGID (group leader)
  - session: PIDTYPE\_SID (session leader)
- *Chaining* is used to handle PID *collisions*

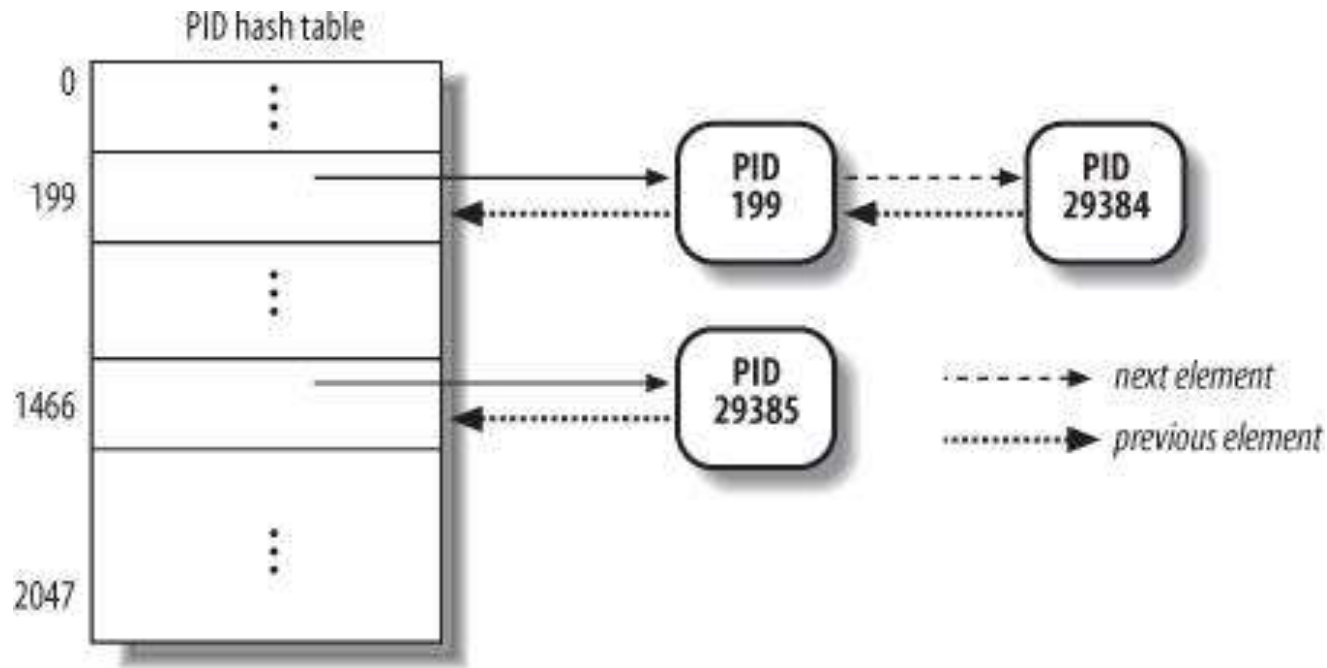


# The pidhash Table

- Size of each pidhash table: dependent on the available memory
- PID is transformed into table index using pid\_hashfn macro
  - #define pid\_hashfn(x) hash\_long((unsigned long)x, pidhash\_shift)
  - unsigned long hash\_long(unsigned long val, unsigned int bits)  
{  
    unsigned long hash = val \* 0x9e370001UL;  
    return hash >> (32-bits);  
}



# A Simple Example PID Hash Table and Chained Lists

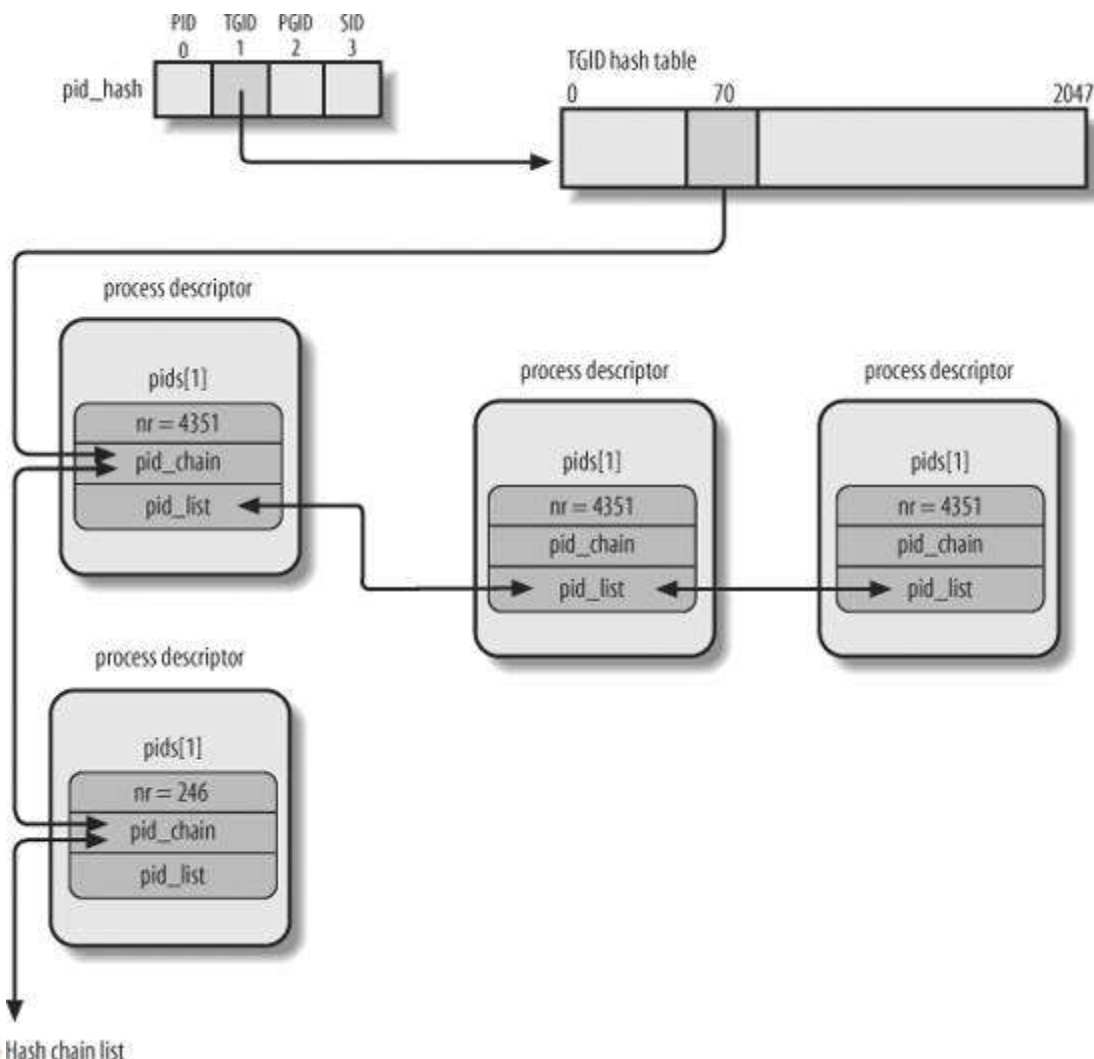


# PID

- *pids* field of the process descriptor: the pid data structures
  - nr: PID number
  - pid\_chain: links to the previous and the next elements in the hash chain list
  - pid\_list: head of the per-PID list (in thread group)



# The PID Hash Tables



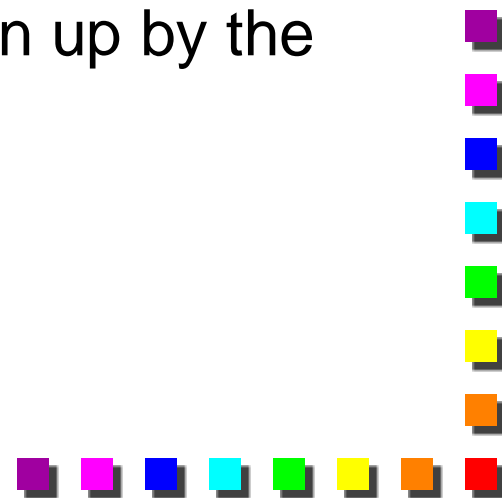
# PID Hash Table Handling Functions and Macros

- `do_each_trask_pid(nr, type, task)`
- `while_each_trask_pid(nr, type, task)`
- `find_trask_by_pid_type(type, nr)`
- `find_trask_by_pid(nr)`
- `attach_pid(task, type, nr)`
- `detach_pid(task, type)`
- `next_thread(task)`



# How Processes are Organized

- Processes in TASK\_STOPPED, EXIT\_ZOMBIE, EXIT\_DEAD: not linked in lists
- Processes in TASK\_INTERRUPTABLE, TASK\_UNINTERRUPTABLE: *wait queues*
- Two kinds of sleeping processes
  - *Exclusive process*
  - *Nonexclusive process*: always woken up by the kernel when the event occurs





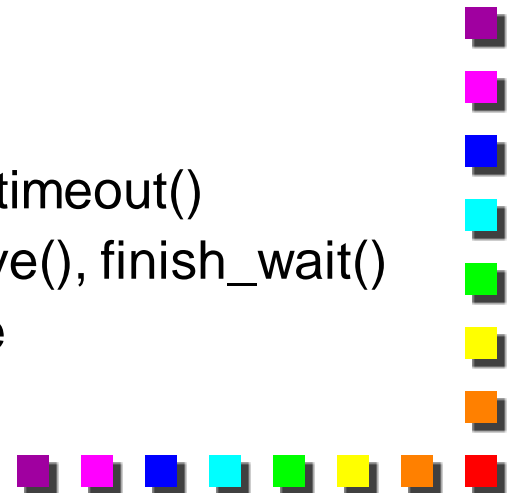
# Wait Queues

- ```
struct __wait_queue_head {  
    spinlock_t lock;  
    struct list_head task_list;  
};  
typedef struct __wait_queue_head  
wait_queue_head_t;
```
- ```
struct __wait_queue {  
    unsigned int flags;  
    struct task_struct * task;  
    wait_queue_func_t func;  
    struct list_head task_list;  
};  
typedef struct __wait_queue wait_queue_t;
```



# Handling Wait Queues (1)

- Wait queue handling functions:
  - `add_wait_queue()`
  - `add_wait_queue_exclusive()`
  - `remove_wait_queue()`
  - `wait_queue_active()`
  - `DECLARE_WAIT_QUEUE_HEAD(name)`
  - `init_waitqueue_head()`
- To wait:
  - `sleep_on()`
  - `interruptible_sleep_on()`
  - `sleep_on_timeout()`, `interruptible_sleep_on_timeout()`
  - `Prepare_to_wait()`, `prepare_to_wait_exclusive()`, `finish_wait()`
  - Macros: `wait_event`, `wait_event_interruptible`



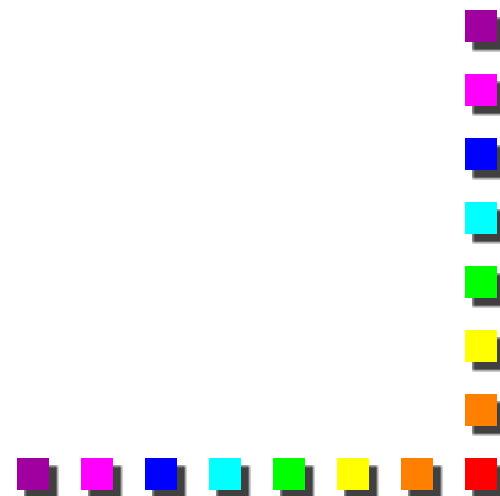
# Handling Wait Queues (2)

- To be woken up:
  - Wake\_up, wake\_up\_nr, wake\_up\_all, wake\_up\_sync, wake\_up\_sync\_nr, wake\_up\_interruptible, wake\_up\_interruptible\_nr, wake\_up\_interruptible\_all, wake\_up\_interruptible\_sync, wake\_up\_interruptible\_sync\_nr



# Process Resource Limits

- RLIMIT\_AS
- RLIMIT\_CORE
- RLIMIT\_CPU
- RLIMIT\_DATA
- RLIMIT\_FSIZE
- RLIMIT\_LOCKS
- RLIMIT\_MEMLOCK
- RLIMIT\_MSGQUEUE
- RLIMIT\_NOFILE
- RLIMIT\_NPROC
- RLIMIT\_RSS
- RLIMIT\_SIGPENDING
- RLIMIT\_STACK



# Process Switch

- Process switch, task switch, *context switch*
  - *Hardware* context switch: a far jmp (in older Linux)
  - *Software* context switch: a sequence of mov instructions
    - It allows better control over the validity of data being loaded
    - The amount of time required is about the same
- Performing the Process Switch
  - Switching the Page Global Directory
  - Switching the Kernel Mode stack and the hardware context



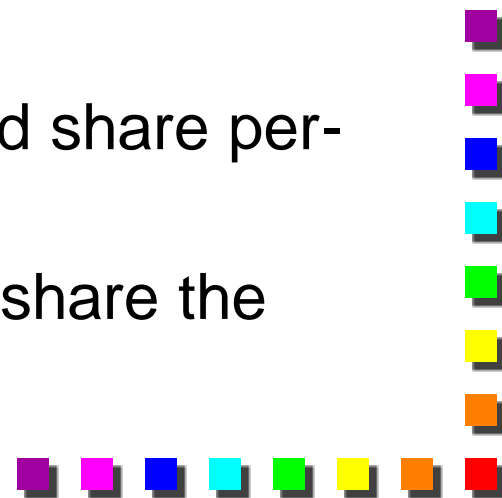
# Task State Segment

- TSS: a specific segment type in x86 architecture to store hardware contexts



# Creating Processes

- In traditional UNIX, resources owned by parent process are duplicated
  - Very slow and inefficient
- Mechanisms to solve this problem
  - **Copy on Write**: parent and child read the same physical pages
  - Lightweight process: parent and child share per-process kernel data structures
  - `vfork()` system call: parent and child share the memory address space



# clone(), fork(), and vfork() System Calls

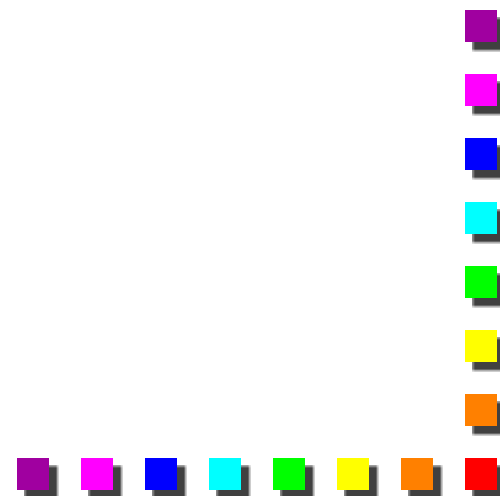
- clone(fn, arg, flags, child\_stack, tls, ptid, ctid):  
creating lightweight process
  - A wrapper function in C library
  - Uses clone() system call
- fork() and vfork() system calls: implemented by clone() with different parameters
- Each invokes *do\_fork()* function





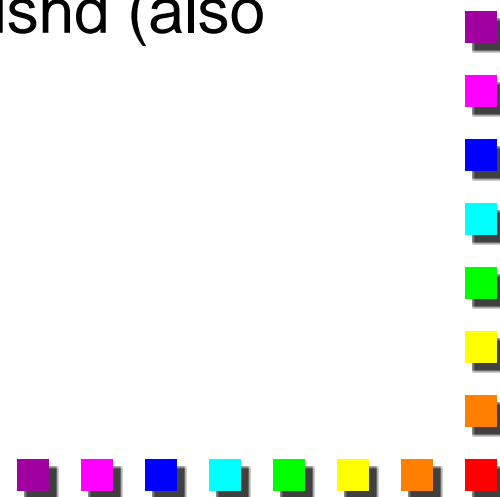
# Kernel Threads

- Kernel threads run only in kernel mode
- They use only linear addresses greater than PAGE\_OFFSET



# Kernel Threads

- `kernel_thread()`: to create a kernel thread
- Example kernel threads
  - Process 0 (swapper process), the ancestor of all processes
  - Process 1 (init process)
  - Others: `keventd`, `kapm`, `kswapd`, `kflushd` (also `bdfush`), `kupdated`, `ksoftirqd`, ...
  -

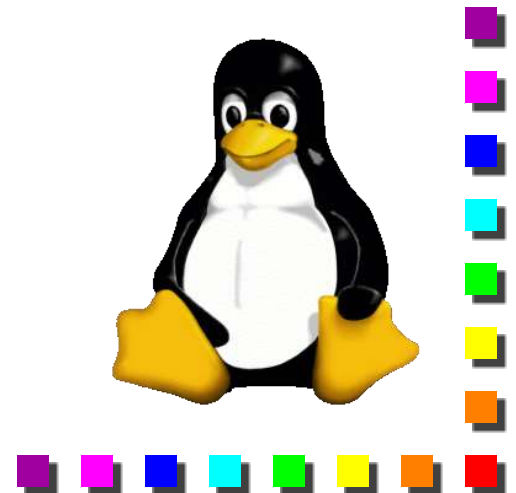


# Destroying Processes

- `exit()` library function
  - Two system calls in Linux 2.6
    - `_exit()` system call
      - Handled by `do_exit()` function
    - `exit_group()` system call
      - By `do_group_exit()` function
- Process removal
  - Releasing the process descriptor of a zombie process by `release_task()`



# Project 2B: Process Management



上海交通大學

# Process: schedule in times

- Add ctx, a new member to task\_struct to record the schedule in times of the process;
  - When a task is scheduled in to run on a CPU, increase ctx of the process;
  - Export ctx under /proc/XXX/ctx;
- More detailed information is shown in **Student Experimental Handbook**.

