

Synchronization

My formulation

- **OS = data structures + synchronization**
- **Synchronization problems make writing OS code challenging**
- **Demand exceptional coding skills**

Race problem

```
long c = 0, c1 = 0, c2 = 0, run = 1;
void *thread1(void *arg) {
    while(run == 1) {
        c++;
        c1++;
    }
}
void *thread2(void *arg) {
    while(run == 1) {
        c++;
        c2++;
    }
}
```

```
int main() {
    pthread_t th1, th2;
    pthread_create(&th1, NULL, thread1,
    NULL);
    pthread_create(&th2, NULL, thread2,
    NULL);
    //fprintf(stdout, "Ending main\n");
    sleep(2);
    run = 0;
    fprintf(stdout, "c = %ld c1+c2 = %ld
c1 = %ld c2 = %ld \n", c, c1+c2, c1, c2);
    fflush(stdout);
}
```

Race problem

- On earlier slide
 - Value of c should be equal to $c1 + c2$, but it is not!
 - Why?
- There is a “race” between thread1 and thread2 for updating the variable c
- thread1 and thread2 may get scheduled in any order and *interrupted* any point in time
- The changes to c are not atomic!
 - What does that mean?

Race problem

- **C++, when converted to assembly code, could be**

```
mov c, r1  
add r1, 1  
mov r1, c
```

- **Now following sequence of instructions is possible among thread1 and thread2**

```
thread1: mov c, r1  
thread2: mov c, r1  
thread1: add r1, 1  
thread1: mov r1, c  
thread2: add r1, 1  
thread2: mov r1, c
```

- **What will be value in c, if initially c was, say 5?**

- It will be 6, when it is expected to be 7. Other variations also possible.

Races: reasons

- **Interruptible kernel**
 - If entry to kernel code does not disable interrupts, then modifications to any kernel data structure can be left incomplete
 - This introduces concurrency
- **Multiprocessor systems**
 - On SMP systems: memory is shared, kernel and process code run on all processors
 - Same variable can be updated parallelly (not concurrently)
- **What about non-interruptible kernel on multiprocessor systems?**
- **What about non-interruptible kernel on uniprocessor systems?**

Critical Section Problem

- Consider system of n processes $\{p_0, p_1, \dots, p_{n-1}\}$
- Each process has critical section segment of code
 - Process may be changing common variables, updating table, writing file, etc
 - When one process in critical section, no other may be in its critical section
- Critical section problem is to design protocol to solve this
- Each process must ask permission to enter critical section in entry section, may follow critical section with exit section, then remainder section
- Especially challenging with preemptive kernels

Critical Section problem

```
do {  
    entry section  
    critical section  
    exit section  
    remainder section  
} while (TRUE);
```

Figure 6.1 General structure of a typical process P .

Expected solution characteristics

- **1. Mutual Exclusion**

- If process P_i is executing in its critical section, then no other processes can be executing in their critical sections

- **2. Progress**

- If no process is executing in its critical section and there exist some processes that wish to enter their critical section, then the selection of the processes that will enter the critical section next cannot be postponed indefinitely

- **3. Bounded Waiting**

- A bound must exist on the number of times that other processes are allowed to enter their critical sections after a process has made a request to enter its critical section and before that request is granted
 - Assume that each process executes at a nonzero speed
 - No assumption concerning relative speed of the n processes

suggested solution - 1

```
int flag = 1;
void *thread1(void *arg) {
    while(run == 1) {
        while(flag == 0)
            ;
        flag = 0;
        c++;
        flag = 1;
        c1++;
    }
}
```

- **What's wrong here?**
- **Assumes that**
while(flag ==) ; flag
= 0
will be atomic

suggested solution - 2

```
int flag = 0;
void *thread1(void *arg) {
    while(run == 1) {
        if(flag)
            c++;
        else
            continue;
        c1++;
        flag = 0;
    }
}
```

```
void *thread2(void *arg) {
    while(run == 1) {
        if(!flag)
            c++;
        else
            continue;
        c2++;
        flag = 1;
    }
}
```

Peterson's solution

- Two process solution
- Assume that the LOAD and STORE instructions are atomic; that is, cannot be interrupted
- The two processes share two variables:
 int turn;
 Boolean flag[2]
- The variable turn indicates whose turn it is to enter the critical section
- The flag array is used to indicate if a process is ready to enter the critical section. $\text{flag}[i] = \text{true}$ implies that process P_i is ready!

Peterson's solution

```
do {  
    flag[i] = TRUE;  
    turn = j;  
    while (flag[j] && turn == j)  
        ;  
    critical section  
    flag[i] = FALSE;  
    remainder section  
} while (TRUE);
```

- Provable that
 - Mutual exclusion is preserved
 - Progress requirement is satisfied
 - Bounded-waiting requirement is met

Hardware solution – the one actually implemented

- **Many systems provide hardware support for critical section code**
- **Uniprocessors – could disable interrupts**
 - Currently running code would execute without preemption
 - Generally too inefficient on multiprocessor systems
 - Operating systems using this not broadly scalable
- **Modern machines provide special atomic hardware instructions**
 - Atomic = non-interruptable
 - Either test memory word and set value
 - Or swap contents of two memory words
 - Basically two operations (read/write) done atomically in hardware

Solution using test-and-set

```
lock = false; //global
```

```
do {  
    while ( TestAndSet (&lock ))  
        ; // do nothing  
    //    critical section  
    lock = FALSE;  
    //    remainder section  
} while (TRUE);
```

Definition:

```
boolean TestAndSet (boolean  
    *target)  
{  
    boolean rv = *target;  
    *target = TRUE;  
    return rv;  
}
```

Solution using swap

```
lock = false; //global
```

```
do {  
    key = true  
    while ( key == true)  
        swap(&lock, &key)  
    //    critical section  
    lock = FALSE;  
    //    remainder section  
} while (TRUE);
```


Spinlock

- A lock implemented to do 'busy-wait'
- Using instructions like T&S or Swap
- As shown on earlier slides

```
spinlock(int *lock){  
    While(test-and-set(lock))  
        ;  
}  
spinunlock(lock *lock) {  
    *lock = false;  
}
```

Bounded wait M.E. with T&S

```
do {  
    waiting[i] = TRUE;  
    key = TRUE;  
    while (waiting[i] && key)  
        key = TestAndSet(&lock);  
    waiting[i] = FALSE;  
    // critical section  
    j = (i + 1) % n;  
    while ((j != i) && !waiting[j])  
        j = (j + 1) % n;  
    if (j == i)  
        lock = FALSE;  
    else  
        waiting[j] = FALSE;  
    // remainder section  
} while (TRUE);
```

Some thumb-rules of spinlocks

- **Never block a process holding a spinlock !**
- **Typical code:**

```
while(condition)
    { Spin-unlock()
      Schedule()
      Spin-lock()
    }
```
- **Hold a spin lock for only a short duration of time**
 - Spinlocks are preferable on multiprocessor systems
 - Cost of context switch is a concern in case of sleep-wait locks
 - Short = < 2 context switches

sleep-locks

- **Spin locks result in busy-wait**
- **CPU cycles wasted by waiting processes/threads**
- **Solution – threads keep waiting for the lock to be available**
 - **Move thread to wait queue**
 - **The thread holding the lock will wake up one of them**

Sleep locks/mutexes

//ignore syntactical issues

typedef struct mutex {

int islocked;

int spinlock;

waitqueue q;

}mutex;

wait(mutex *m) {

spinlock(m->spinlock);

while(m->islocked)

Block(m, m->spinlock)

lk->islocked = 1;

spinunlock(m->spinlock);

}

Block(mutex *m, spinlock *sl) {

spinunlock(sl);

currprocess->state = WAITING

move current process to m->q

Sched();

spinlock(sl);

}

release(mutex *m) {

spinlock(m->spinlock);

m->islocked = 0;

Some process in m->queue
=RUNNABLE;

spinunlock(m->spinlock);

}

Locks in xv6 code

struct spinlock

// Mutual exclusion lock.

struct spinlock {

uint locked; // Is the lock held?

// For debugging:

char *name; // Name of lock.

struct cpu *cpu; // The cpu holding the lock.

uint pcs[10]; // The call stack (an array of program counters)
 // that locked the lock.

};

spinlocks in xv6 code

```
struct {  
    struct spinlock lock;  
    struct buf buf[NBUF];  
    struct buf head;  
} bcache;  
  
struct {  
    struct spinlock lock;  
    struct file file[NFILE];  
} ftable;  
  
struct {  
    struct spinlock lock;  
    struct inode inode[NINODE];  
} icache;  
  
struct sleeplock {  
    uint locked;    // Is the lock held?  
    struct spinlock lk;
```

```
static struct spinlock idelock;  
  
struct {  
    struct spinlock lock;  
    int use_lock;  
    struct run *freelist;  
} kmem;  
  
struct log {  
    struct spinlock lock;  
    ...}  
  
struct pipe {  
    struct spinlock lock;  
    ...}  
  
struct {  
    struct spinlock lock;  
    struct proc proc[NPROC];  
} ptable;  
  
struct spinlock tickslock;
```



```

static inline uint
xchg(volatile uint *addr, uint newval)
{
    uint result;

    // The + in "+m" denotes a read-modify-
    write operand.
    asm volatile("lock; xchgl %0, %1" :
        "+m" (*addr), "=a" (result) :
        "1" (newval) :
        "cc");
    return result;
}

struct spinlock {
    uint locked;    // Is the lock held?

    // For debugging:
    char *name;     // Name of lock.
    struct cpu *cpu; // The cpu holding the
    lock.

    uint pcs[10];   // The call stack (an array
    of program counters) that locked the lock.
};

```

Spinlock on xv6

```

void acquire(struct spinlock *lk)
{
    pushcli(); // disable interrupts to
    avoid deadlock.

    // The xchg is atomic.
    while(xchg(&lk->locked, 1) != 0)
        ;

    //extra debugging code
}

void release(struct spinlock *lk)
{
    //extra debugging code
    asm volatile("movl $0, %0" :
        "+m" (lk->locked) : );
    popcli();
}

```

```
Void acquire(struct spinlock *lk)
```

```
{  
    pushcli(); // disable interrupts to avoid deadlock.  
    if(holding(lk))  
        panic("acquire");  
    .....
```

```
void pushcli(void)
```

```
{  
    int eflags;  
  
    eflags = readeflags();  
    cli();  
    if(mycpu()->ncli == 0)  
        mycpu()->intena = eflags & FL_IF;  
    mycpu()->ncli += 1;  
}
```

```
static inline uint
```

```
readeflags(void)
```

```
{  
    uint eflags;  
    asm volatile("pushfl; popl %0" : "=r" (eflags));  
    return eflags;  
}
```

spinlocks

- **Pushcli() - disable interrupts on that processor**
- **One after another many acquire() can be called on different spinlocks**
 - **Keep a count of them in mycpu()->ncli**

```

void
release(struct spinlock *lk)
{
...
    asm volatile("movl $0, %0" : "+m" (lk-
>locked) : );
    popcli();
}
.
Void popcli(void)
{
    if(readeflags() & FL_IF)
        panic("popcli - interruptible");
    if(--mycpu()->ncli < 0)
        panic("popcli");
    if(mycpu()->ncli == 0 && mycpu()->intena)
        sti();
}

```

spinlocks

- **Popcli()**
 - Restore interrupts if last popcli() call restores ncli to 0 & interrupts were enabled before pushcli() was called

spinlocks

- **Always disable interrupts while acquiring spinlock**
 - Suppose **iderw** held the **idelock** and then got interrupted to run **ideintr**.
 - **Ideintr** would try to lock **idelock**, see it was held, and wait for it to be released.
 - In this situation, **idelock** will never be released
 - Deadlock
- **General OS rule: if a spin-lock is used by an interrupt handler, a processor must never hold that lock with interrupts enabled**
- **Xv6 rule: when a processor enters a spin-lock critical section, xv6 always ensures interrupts are disabled on that processor.**

sleeplocks

- **Sleeplocks don't spin. They move a process to a wait-queue if the lock can't be acquired**
- **XV6 approach to “wait-queues”**
 - Any memory address serves as a “wait channel”
 - The sleep() and wakeup() functions just use that address as a ‘condition’
 - There are no per condition process queues! Just one global queue of processes used for scheduling, sleep, wakeup etc. --> Linear search everytime !
 - **costly, but simple**

void

sleep(void *chan, struct spinlock *lk)

{

struct proc *p = myproc();

....

if(lk != &ptable.lock){

acquire(&ptable.lock);

release(lk);

}

p->chan = chan;

p->state = SLEEPING;

sched();

// Reacquire original lock.

if(lk != &ptable.lock){

release(&ptable.lock);

acquire(lk);

}

sleep()

- At call must hold lock on the resource on which you are going to sleep
- since you are going to change p-> values & call sched(), hold ptable.lock if not held
- p->chan = given address remembers on which condition the process is waiting
- call to sched() blocks the process

Calls to sleep() : examples of “chan” (output from cscope)

0 console.c
consoleread 251
sleep(&input.r, &cons.lock);

2 ide.c iderw
169 sleep(b, &idelock);

3 log.c begin_op
131 sleep(&log, &log.lock);

6 pipe.c piperead
111 sleep(&p->nread, &p->lock);

7 proc.c wait
317 sleep(curproc, &ptable.lock);

8 sleeplock.c
acquiresleep 28
sleep(lk, &lk->lk);

9 sysproc.c
sys_sleep 74
sleep(&ticks, &tickslock);

```

void wakeup(void *chan)
{
    acquire(&ptable.lock);
    wakeup1(chan);
    release(&ptable.lock);
}

static void wakeup1(void *chan)
{
    struct proc *p;

    for(p = ptable.proc; p <
        &ptable.proc[NPROC]; p++)
        if(p->state == SLEEPING &&
            p->chan == chan)
            p->state = RUNNABLE;
}

```

Wakeup()

- **Acquire ptable.lock since you are going to change ptable and p->values**
- **just linear search in process table for a process where p->chan is given address**
- **Make it runnable**

sleeplock

// Long-term locks for processes

struct sleeplock {

uint locked; // Is the lock held?

struct spinlock lk; // spinlock protecting this sleep lock

// For debugging:

char *name; // Name of lock.

int pid; // Process holding lock

};

Sleeplock acquire and release

```
void
acquiresleep(struct sleeplock *lk)
{
    acquire(&lk->lk);
    while (lk->locked) {
        /* Abhijit: interrupts are not disabled in
        sleep !*/
        sleep(lk, &lk->lk);
    }
    lk->locked = 1;
    lk->pid = myproc()->pid;
    release(&lk->lk);
}
```

```
void
releasesleep(struct sleeplock
*lk)
{
    acquire(&lk->lk);
    lk->locked = 0;
    lk->pid = 0;
    wakeup(lk);
    release(&lk->lk);
}
```

Where are sleeplocks used?

- **struct buf**
 - waiting for I/O on this buffer
- **struct inode**
 - waiting for I/o to this inode
- **Just two !**

Sleeplocks issues

- **sleep-locks support yielding the processor during their critical sections.**
- **This property poses a design challenge:**
 - if thread T1 holds lock L1 and has yielded the processor (waiting for some other condition),
 - and thread T2 wishes to acquire L1,
 - we have to ensure that T1 can execute
 - while T2 is waiting so that T1 can release L1.
 - T2 can't use the spin-lock acquire function here: it spins with interrupts turned off, and that would prevent T1 from running.
- **To avoid this deadlock, the sleep-lock acquire routine (called `acquiresleep`) yields the processor while waiting, and does not disable interrupts.**

Sleep-locks leave interrupts enabled, they cannot be used in interrupt handlers.

More needs of synchronization

- Not only critical section problems
- Run processes in a particular order
- Allow multiple processes read access, but only one process write access
- Etc.

Semaphore

- Synchronization tool that does not require busy waiting
 - Semaphore S – integer variable
 - Two standard operations modify S: wait() and signal()
 - Originally called P() and V()
 - Less complicated
- Can only be accessed via two indivisible (atomic) operations

```
wait (S) {  
    while S <= 0  
        ; // no-op  
    S--;  
}  
signal (S) {  
    S++;  
}
```

--> Note this is Signal() on a semaphore, different froms signal system call

Semaphore for synchronization

- **Counting** semaphore – integer value can range over an unrestricted domain
- **Binary** semaphore – integer value can range only between 0 and 1; can be simpler to implement

Also known as **mutex locks**

- Can implement a counting semaphore **S** as a binary semaphore
- Provides mutual exclusion

```
Semaphore mutex; // initialized to 1
do {
    wait (mutex);
        // Critical Section
    signal (mutex);
// remainder section
} while (TRUE)
```

Semaphore implementation

```
Wait(sem *s) {  
    while(s <=0)  
        block(); // could be ";"  
    s--;  
}  
  
signal(sem *s) {  
    s++;  
}
```

- Left side – expected behaviour
- Both the wait and signal should be atomic.
- This is the semantics of the semaphore.

Semaphore implementation? - 1

```
struct semaphore {  
    int val;  
    spinlock lk;  
};  
sem_init(semaphore *s, int initval) {  
    s->val = initval;  
    s->sl = 0;  
}  
wait(semaphore *s) {  
    spinlock(&(s->sl));  
    while(s->val <= 0)  
        ;  
    (s->val)--;  
    spinunlock(&(s->sl));  
}
```

```
signal(semaphore *s) {  
    spinlock(&(s->sl));  
    (s->val)++;  
    spinunlock(&(s->sl));  
}
```

- suppose 2 processes trying wait.

val = 1;

Th1: spinlock Th2: spinlock-waits

Th1: while -> false, val-- => 0; spinunlock;

Th2: spinlock success; while() -> true, loops;

Th1: is done with critical section, it calls signal. it calls spinlock() -> wait.

Who is holding spinlock-> Th2. It is waiting for val > 0. Who can set value > 0, ans: Th1, and Th1 is waiting for spinlock which is held by Th2.

circular wait. Deadlock.

None of them will proceed.

Semaphore implementation? - 2

```
struct semaphore {  
    int val;  
    spinlock lk;  
};  
  
sem_init(semaphore *s, int initval) {  
    s->val = initval;  
    s->sl = 0;  
}  
  
signal(semaphore *s) {  
    spinlock(&(s->sl));  
    (s->val)++;  
    spinunlock(&(s->sl));  
}
```

```
wait(semaphore *s) {  
    spinlock(&(s->sl));  
    while(s->val <= 0) {  
        spinunlock(&(s->sl));  
        spinlock(&(s->sl));  
    }  
    (s->val)--;  
    spinunlock(&(s->sl));  
}
```

**Problem: race in spinlock of while loop and signal's spinlock.
Bounded wait not guaranteed.**

Spinlocks are not good for a long wait.

Semaphore implementation? - 3, idea

```
struct semaphore {
    int val;
    spinlock lk;
};
sem_init(semaphore *s, int initval) {
    s->val = initval;
    s->sl = 0;
}
block() {
    put this current process on wait-q;
    schedule();
}
```

```
wait(semaphore *s) {
    spinlock(&(s->sl));
    while(s->val <= 0) {
        Block();
    }
    (s->val)--;
    spinunlock(&(s->sl));
}
signal(semaphore *s) {
    spinlock(&(s->sl));
    (s->val)++;
    spinunlock(&(s->sl));
}
```

Semaphore implementation? - 3a

```
struct semaphore {  
    int val;  
    spinlock lk;  
    list l;  
};  
sem_init(semaphore *s, int initval) {  
    s->val = initval;  
    s->sl = 0;  
}  
block(semaphore *s) {  
    listappend(s->l, current);  
    schedule();  
}
```

problem is that block() will be called without holding the spinlock and the access to the list is not protected.

Note that - so far we have ignored changes to signal()

```
wait(semaphore *s) {  
    spinlock(&(s->sl));  
    while(s->val <= 0) {  
        spinunlock(&(s->sl));  
        block(s);  
    }  
    (s->val)--;  
    spinunlock(&(s->sl));  
}  
signal(semaphore *s) {  
    spinlock(&(s->sl));  
    (s->val)++;  
    spinunlock(&(s->sl));  
}
```

Semaphore implementation? - 3b

```
struct semaphore {  
    int val;  
    spinlock lk;  
    list l;  
};  
sem_init(semaphore *s, int initval) {  
    s->val = initval;  
    s->sl = 0;  
}  
block(semaphore *s) {  
    listappend(s->l, current);  
    spinunlock(&(s->sl));  
    schedule();  
}
```

```
wait(semaphore *s) {  
    spinlock(&(s->sl));  
    while(s->val <= 0) {  
        block(s);  
    }  
    (s->val)--;  
    spinunlock(&(s->sl));  
}  
signal(semaphore *s) {  
    spinlock(&(s->sl));  
    (s->val)++;  
    x = dequeue(s->sl) and enqueue(readyq, x);  
    spinunlock(&(s->sl));  
}  
Problem: after a blocked process comes out  
of the block, it does not hold the spinlock and  
it's going to change the s->sl;
```

Semaphore implementation? - 3c

```
struct semaphore {  
    int val;  
    spinlock lk;  
    list l;  
};  
  
sem_init(semaphore *s, int initval) {  
    s->val = initval;  
    s->sl = 0;  
}  
  
block(semaphore *s) {  
    listappend(s->l, current);  
    spinunlock(&(s->sl));  
    schedule();  
}
```

```
wait(semaphore *s) {  
    spinlock(&(s->sl)); // A  
    while(s->val <= 0) {  
        block(s);  
        spinlock(&(s->sl)); // B  
    }  
    (s->val)--;  
    spinunlock(&(s->sl));  
}  
  
signal(semaphore *s) {  
    spinlock(&(s->sl));  
    (s->val)++;  
    x = dequeue(s->sl) and enqueue(readyq, x);  
    spinunlock(&(s->sl));  
}
```

Question: there is race between A and B. Can we guarantee bounded wait ?

Semaphore Implementation

- Must guarantee that no two processes can execute wait () and signal () on the same semaphore at the same time
- Thus, implementation becomes the critical section problem where the wait and signal code are placed in the critical section
 - Could now have busy waiting in critical section implementation
 - But implementation code is short
 - Little busy waiting if critical section rarely occupied
- Note that applications may spend lots of time in critical sections and therefore this is not a good solution

Semaphore in Linux

```
struct semaphore {  
    raw_spinlock_t    lock;  
    unsigned int      count;  
    struct list_head   wait_list;  
};  
  
static ninline void __sched  
__down(struct semaphore *sem)  
{  
    __down_common(sem,  
TASK_UNINTERRUPTIBLE,  
MAX_SCHEDULE_TIMEOUT);  
}
```

```
void down(struct semaphore *sem)  
{  
    unsigned long flags;  
  
    raw_spin_lock_irqsave(&sem->lock, flags);  
    if (likely(sem->count > 0))  
        sem->count--;  
    else  
        __down(sem);  
    raw_spin_unlock_irqrestore(&sem->lock,  
flags);  
}
```


Semaphore in Linux

```
static inline int __sched
__down_common(struct semaphore
*sem, long state, long timeout)
{
    struct task_struct *task = current;
    struct semaphore_waiter waiter;
    list_add_tail(&waiter.list, &sem-
>wait_list);
    waiter.task = task;
    waiter.up = false;
```

```
    for (;;) {
        if (signal_pending_state(state, task))
            goto interrupted;
        if (unlikely(timeout <= 0))
            goto timed_out;
        __set_task_state(task, state);
        raw_spin_unlock_irq(&sem->lock);
        timeout = schedule_timeout(timeout);
        raw_spin_lock_irq(&sem->lock);
        if (waiter.up)
            return 0;
    }
    ....
}
```

Different uses of semaphores

For mutual exclusion

/*During inialization*/

semaphore sem;

initsem (&sem, 1);

/* On each use*/

P (&sem);

Use resource;

V (&sem);

Event-wait

/* During initialization */

semaphore event;

initsem (&event, 0); /* probably at boot time */

/* Code executed by thread that must wait on event */

P (&event); /* Blocks if event has not occurred */

/* Event has occurred */

V (&event); /* So that another thread may wake up */

/* Continue processing */

/* Code executed by another thread when event occurs */

V (&event); /* Wake up one thread */

Control countable resources

/* During initialization */

semaphore counter;

initsem (&counter, resourceCount);

/* Code executed to use the resource */

P (&counter); /* Blocks until resource is available */

Use resource; /* Guaranteed to be available now */

V (&counter); /* Release the resource */

Drawbacks of semaphores

- **Need to be implemented using lower level primitives like spinlocks**
- **Context-switch is involved in blocking and signaling – time consuming**
- **Can not be used for a short critical section**

Deadlocks

Deadlock

- two or more processes are waiting indefinitely for an event that can be caused by only one of the waiting processes
- Let S and Q be two semaphores initialized to 1

P0

wait (S);

wait (Q);

▪ ▪

▪ ▪

▪ ▪

signal (S);

signal (Q);

P1

wait (Q);

wait (S);

signal (Q);

signal (S);

Example of deadlock

- Let's see the pthreads program : deadlock.c
- Same programme as on earlier slide, but with `pthread_mutex_lock()`;

Non-deadlock, but similar situations

- **Starvation – indefinite blocking**
 - A process may never be removed from the semaphore queue in which it is suspended
- **Priority Inversion**
 - Scheduling problem when lower-priority process holds a lock needed by higher-priority process (so it can not pre-empt lower priority process), and a medium priority process (that does not need the lock) pre-empts lower priority task, denying turn to higher priority task
 - Solved via priority-inheritance protocol : temporarily enhance priority of lower priority task to highest

Livelock

- **Similar to deadlock, but processes keep doing 'useless work'**
- **E.g. two people meet in a corridor opposite each other**
 - **Both move to left at same time**
 - **Then both move to right at same time**
 - **Keep Repeating!**
- **No process able to progress, but each doing 'some work' (not sleeping/waiting), state keeps changing**

Livelock example

```
#include <stdio.h>
#include <pthread.h>
struct person {
    int otherid;
    int otherHungry;
    int myid;
};
int main() {
    pthread_t th1, th2;
    struct person one, two;
    one.otherid = 2; one.myid = 1;
    two.otherid = 1; two.myid = 2;
    one.otherHungry = two.otherHungry = 1;
    pthread_create(&th1, NULL, eat, &one);
    pthread_create(&th2, NULL, eat, &two);
    printf("Main: Waiting for threads to get over\n");
    pthread_join(th1, NULL);
    pthread_join(th2, NULL);
    return 0;
}
```

```
/* thread two runs in this function */
int spoonWith = 1;
void *eat(void *param)
{
    int eaten = 0;
    struct person person= *(struct person *)param;
    while (!eaten) {
        if(spoonWith == person.myid)
            printf("%d going to eat\n", person.myid);
        else
            continue;
        if(person.otherHungry) {
            printf("You eat %d\n", person.otherid);
            spoonWith = person.otherid;
            continue;
        }
        printf("%d is eating\n", person.myid);
        break;
    }
}
```

More on deadlocks

- Under which conditions they can occur?
- How can deadlocks be avoided/prevented?
- How can a system recover if there is a deadlock ?

System model for understanding deadlocks

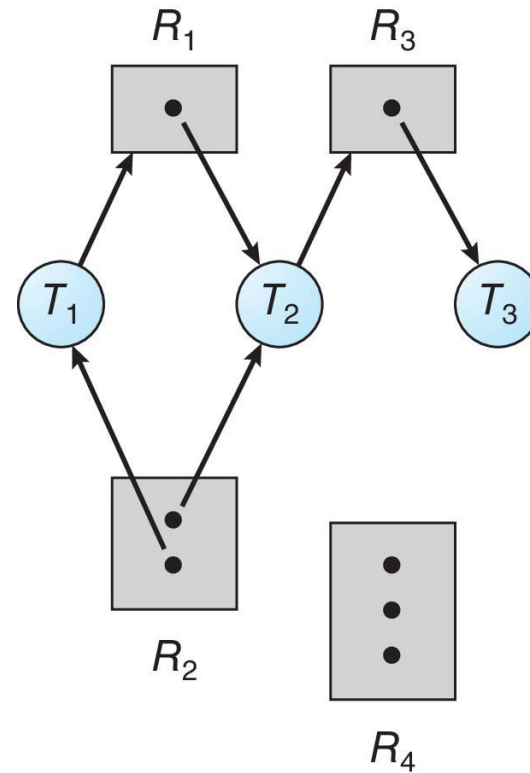
- **System consists of resources**
- **Resource types R_1, R_2, \dots, R_m**
 - CPU cycles, memory space, I/O devices
 - Resource: Most typically a lock, synchronization primitive
- **Each resource type R_i has W_i instances.**
- **Each process utilizes a resource as follows:**
 - request
 - use
 - release

Deadlock characterisation

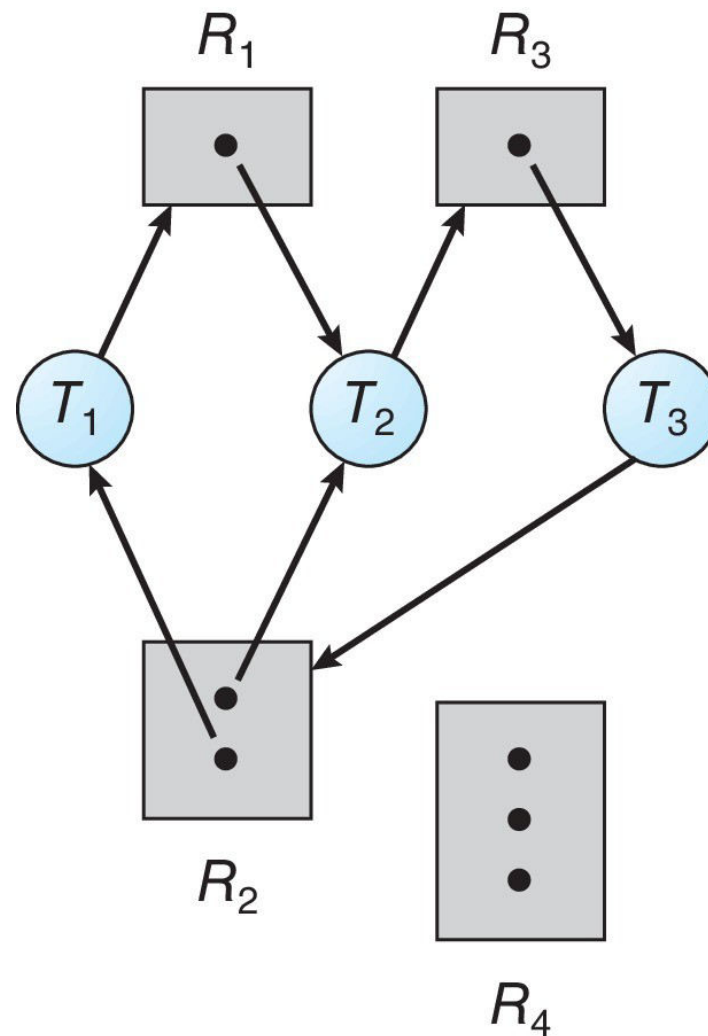
- **Deadlock is possible only if ALL of these conditions are TRUE at the same time**
 - **Mutual exclusion:** only one process at a time can use a resource
 - **Hold and wait:** a process holding at least one resource is waiting to acquire additional resources held by other processes
 - **No preemption:** a resource can be released only voluntarily by the process holding it, after that process has completed its task
 - **Circular wait:** there exists a set $\{P_0, P_1, \dots, P_n\}$ of waiting processes such that P_0 is waiting for a resource that is held by P_1 , P_1 is waiting for a resource that is held by P_2 , ..., P_{n-1} is waiting for a resource that is held by P_n , and P_n is waiting for a resource that is held by P_0 .

Resource Allocation Graph Example

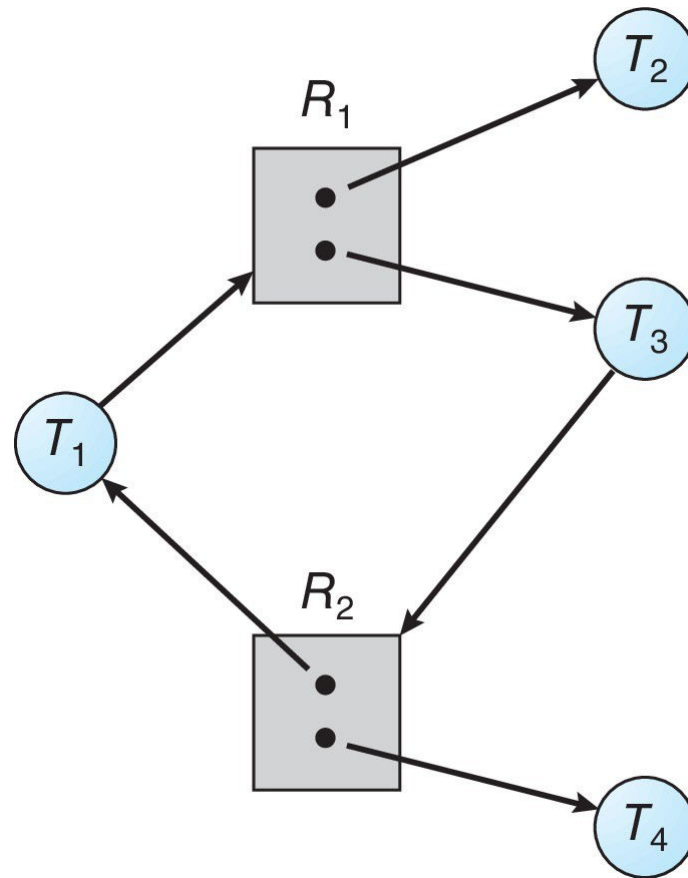
- One instance of R1
- Two instances of R2
- One instance of R3
- Three instance of R4
- T1 holds one instance of R2 and is waiting for an instance of R1
- T2 holds one instance of R1, one instance of R2, and is waiting for an instance of R3
- T3 is holds one instance of R3



Resource Allocation Graph with a Deadlock



Graph with a Cycle But no Deadlock



Basic Facts

- **If graph contains no cycles -> no deadlock**
- **If graph contains a cycle :**
 - **if only one instance per resource type, then deadlock**
 - **if several instances per resource type, possibility of deadlock**

Methods for Handling Deadlocks

- **Ensure that the system will never enter a deadlock state:**
 - 1) Deadlock prevention**
 - 2) Deadlock avoidance**
 - 3) Allow the system to enter a deadlock state and then recover**
 - 4) Ignore the problem and pretend that deadlocks never occur in the system.**

(1) Deadlock Prevention

- Invalidate one of the four necessary conditions for deadlock:
- **Mutual Exclusion** – not required for sharable resources (e.g., read-only files); must hold for non-sharable resources
- **Hold and Wait** – must guarantee that whenever a process requests a resource, it does not hold any other resources
 - Require process to request and be allocated all its resources before it begins execution, or allow process to request resources only when the process has none allocated to it.
 - Low resource utilization; starvation possible

(1) Deadlock Prevention (Cont.)

- **No Preemption:**

- If a process that is holding some resources requests another resource that cannot be immediately allocated to it, then all resources currently being held are released
- Preempted resources are added to the list of resources for which the process is waiting
- Process will be restarted only when it can regain its old resources, as well as the new ones that it is requesting

- **Circular Wait:**

- Impose a total ordering of all resource types, and require that each process requests resources in an increasing order of enumeration

(1) Deadlock prevention: Circular Wait

- Invalidating the circular wait condition is most common.
- Simply assign each resource (i.e., mutex locks) a unique number.
- Resources must be acquired in order.
- If:
 first_mutex = 1
 second_mutex = 5
code for thread_two could not be written as follows:

```
/* thread_one runs in this function */
void *do_work_one(void *param)
{
    pthread_mutex_lock(&first_mutex);
    pthread_mutex_lock(&second_mutex);
    /**
     * Do some work
     */
    pthread_mutex_unlock(&second_mutex);
    pthread_mutex_unlock(&first_mutex);

    pthread_exit(0);
}

/* thread_two runs in this function */
void *do_work_two(void *param)
{
    pthread_mutex_lock(&second_mutex);
    pthread_mutex_lock(&first_mutex);
    /**
     * Do some work
     */
    pthread_mutex_unlock(&first_mutex);
    pthread_mutex_unlock(&second_mutex);

    pthread_exit(0);
}
```

(1) Preventing deadlock: cyclic wait

- **Locking hierarchy : Highly preferred technique in kernels**
 - **Decide an ordering among all 'locks'**
 - **Ensure that on ALL code paths in the kernel, the locks are obtained in the decided order!**
 - **Poses coding challenges!**
 - **A key differentiating factor in kernels**
 - **Do not look at only the current lock being taken, look at all the locks the code may be holding at any given point in code!**

(1) Prevention in Xv6: Lock Ordering

- **lock on the directory, a lock on the new file's inode, a lock on a disk block buffer, idelock, and ptable.lock.**

(2) Deadlock avoidance

- **Requires that the system has some additional a priori information available**
 - **Simplest and most useful model requires that each process declare the maximum number of resources of each type that it may need**
 - **The deadlock-avoidance algorithm dynamically examines the resource-allocation state to ensure that there can never be a circular-wait condition**
 - **Resource-allocation state is defined by the number of available and allocated resources, and the maximum demands of the processes**

(2) Deadlock avoidance

- **Please see: concept of safe states, unsafe states, Banker's algorithm**

(3) Deadlock detection and recovery

- How to detect a deadlock in the system?
- The Resource-Allocation Graph is a graph. Need an algorithm to detect cycle in a graph.
- How to recover?
 - Abort all processes or abort one by one?
 - Which processes to abort?
 - Priority ?
 - Time spent since forked()?
 - Resources used?
 - Resources needed?
 - Interactive or not?
 - How many need to be terminated?

“Condition” Synchronization Tool

What is condition variable?

- A variable with a sleep queue
- Threads can sleep on it, and wake-up all remaining

```
Struct condition {  
    Proc *next  
    Proc *prev  
    Spinlock *lock  
}
```

Different variables of this type can be used as different
'conditions'

Code for condition variables

```
//Spinlock s is held before calling wait
void wait (condition *c, spinlock_t *s)
(
    spin_lock (&c->listLock);
    add self to the linked list;
    spin_unlock (&c->listLock);
    spin_unlock (s); /* release
    spinlock before blocking */
    swtch(); /* perform context switch */
    /* When we return from swtch, the
    event has occurred */
    spin_lock (s); /* acquire the spin
    lock again */
    return;
)
```

```
void do_signal (condition *c)
/*Wakeup one thread waiting on thecondition*/
{
    spin_lock (&c->listLock);
    remove one thread from linked list, if it is nonempty;
    spin_unlock (&c->listLock);
    if a thread was removed from the list, make it
        runnable;
    return;
}
void do broadcast (condition *c)
/*Wakeup all threads waiting on the condition*/
{
    spin_lock (&c->listLock);
    while (linked list is nonempty) {
        remove a thread from linked list;
        make it runnable;
    }
    spin_unlock (&c->listLock);
}
```

Semaphore implementation using condition variables?

- Is this possible?
- Can we try it?

```
typedef struct semaphore {  
    //something  
    condition c;  
}semaphore;
```

- Now write code for semaphore P() and V()

Classical Synchronization Problems

Bounded-Buffer Problem

- **Producer and consumer processes**
 - N buffers, each can hold one item
- **Producer produces 'items' to be consumed by consumer , in the bounded buffer**
- **Consumer should wait if there are no items**
- **Producer should wait if the 'bounded buffer' is full**

Bounded-Buffer Problem: solution with semaphores

- **Semaphore mutex initialized to the value 1**
- **Semaphore full initialized to the value 0**
- **Semaphore empty initialized to the value N**

Bounded-buffer problem

The structure of the producer process

```
do {  
    // produce an item in nextp  
    wait (empty);  
    wait (mutex);  
    // add the item to the buffer  
    signal (mutex);  
    signal (full);  
} while (TRUE);
```

The structure of the Consumer process

```
do {  
    wait (full);  
    wait (mutex);  
    // remove an item from  
    // buffer to nextc  
    signal (mutex);  
    signal (empty);  
    // consume item in nextc  
} while (TRUE);
```

Bounded buffer problem

- **Example : pipe()**
- **Let's see code of pipe in xv6 – a solution using sleeplocks**

Readers-Writers problem

- **A data set is shared among a number of concurrent processes**
 - Readers – only read the data set; they do not perform any updates
 - Writers – can both read and write
- **Problem – allow multiple readers to read at the same time**
 - Only one single writer can access the shared data at the same time
- **Several variations of how readers and writers are treated – all involve priorities**
- **Shared Data**
 - Data set
 - Semaphore mutex initialized to 1
 - Semaphore wrt initialized to 1
 - Integer readcount initialized to 0

The structure of a writer process

```
do {  
    wait (wrt) ;  
    // writing is performed  
    signal (wrt) ;  
} while (TRUE);
```

Readers-Writers problem

The structure of a reader process

```
do {  
  
    wait (mutex) ;  
    readcount ++ ;  
    if (readcount == 1)  
        wait (wrt) ;  
    signal (mutex)  
    // reading is performed  
    wait (mutex) ;  
    readcount - - ;  
    if (readcount == 0)  
        signal (wrt) ;  
    signal (mutex) ;  
  
} while (TRUE);
```

Readers-Writers Problem

Variations

- **First variation – no reader kept waiting unless writer has permission to use shared object**
- **Second variation – once writer is ready, it performs write asap**
- **Both may have starvation leading to even more variations**
- **Problem is solved on some systems by kernel providing reader-writer locks**

Reader-write lock

- **A lock with following operations on it**
 - Lockshared()
 - Unlockshared()
 - LockExcl()
 - UnlockExcl()
- **Possible additions**
 - Downgrade() -> from excl to shared
 - Upgrade() -> from shared to excl

Code for reader-writer locks

```
struct rwlock {  
    int nActive; /* num of active  
    readers, or -1 if a writer is  
    active */  
  
    int nPendingReads;  
    int nPendingWrites;  
    spinlock_t sl;  
    condition canRead;  
    condition canWrite;  
};
```

```
void lockShared (struct rwlock *r)  
{  
    spin_lock (&r->sl);  
    r->nPendingReads++;  
    if (r->nPendingWrites > 0)  
        wait (&r->canRead, &r->sl); /*don'tstarve  
        writers */  
    while (r->nActive < 0) /* someone has  
        exclusive lock */  
        wait (&r->canRead, &r->sl);  
    r->nActive++;  
    r->nPendingReads--;  
    spin_unlock (&r->sl);  
}
```

Code for reader-writer locks

```
void unlockShared (struct rwlock
*r)
{
    spin_lock (&r->sl);
    r->nActive--;
    if (r->nActive == 0) {
        spin_unlock (&r->sl);
        do signal (&r->canWrite);
    } else
        spin_unlock (&r->M);
}
```

```
void lockExclusive (struct rwlock
*r)
(
    spin_lock (&r->sl);
    r->nPendingWrites++;
    while (r->nActive)
        wait (&r->canWrite, &r->sl);
    r->nPendingWrites--;
    r->nActive = -1;
    spin_unlock (&r->sl);
}
```

Code for reader-writer locks

```
void unlockExclusive (struct rwlock *r){
    boolean t wakeReaders;
    spin_lock (&r->sl);
    r->nActive = 0;
    wakeReaders = (r->nPendingReads != 0);
    spin_unlock (&r->sl);
    if (wakeReaders)
        do broadcast (&r->canRead); /* wake
allreaders */
    else
        do_signal (&r->canWrite);
        /*wakeasinglewrrir */
}
```

**Try writing code for
downgrade and
upgrade**

**Try writing a reader-
writer lock using
semaphores!**

Dining-Philosophers Problem

- Philosophers spend their lives thinking and eating
- Don't interact with their neighbors, occasionally try to pick up 2 chopsticks (one at a time) to eat from bowl
 - Need both to eat, then release both when done
- In the case of 5 philosophers
 - Shared data
 - Bowl of rice (data set)
 - Semaphore chopstick [5] initialized to 1



Dining philosophers: One solution

The structure of Philosopher i:

```
do {  
    wait ( chopstick[i] );  
    wait ( chopStick[ (i + 1) % 5] );  
    // eat  
    signal ( chopstick[i] );  
    signal ( chopstick[ (i + 1) % 5] );  
    // think  
} while (TRUE);
```

What is the problem with this algorithm?

Dining philosophers: Possible approaches

- **Allow at most four philosophers to be sitting simultaneously at the table.**
- **Allow a philosopher to pick up her chopsticks only if both chopsticks are available**
 - to do this, she must pick them up in a critical section
- **Use an asymmetric solution**
 - that is, an odd-numbered philosopher picks up first her left chopstick and then her right chopstick
 - whereas an even-numbered philosopher picks up her right chopstick and then her left chopstick.

Other solutions to dining philosopher's problem

- Using higher level synchronization primitives like 'monitors'
-

Practical Problems

Lost Wakeup problem

- The sleep/wakeup mechanism does not function correctly on a multiprocessor.
- Consider a potential race:
 - Thread T1 has locked a resource R1.
 - Thread T2, running on another processor, tries to acquire the resource, and finds it locked.
 - T2 calls sleep() to wait for the resource.
 - Between the time T2 finds the resource locked and the time it calls sleep (), T1 frees the resource and proceeds to wake up all threads blocked on it.
 - Since T2 has not yet been put on the sleep queue, it will miss the wakeup.
 - The end result is that the resource is not locked, but T2 is blocked waiting for it to be unlocked.
 - If no one else tries to access the resource, T2 could block indefinitely.
 - This is known as the lost wakeup problem,
- Requires some mechanism to combine the test for the resource and the call to sleep () into a single atomic operation.

Lost Wakeup problem

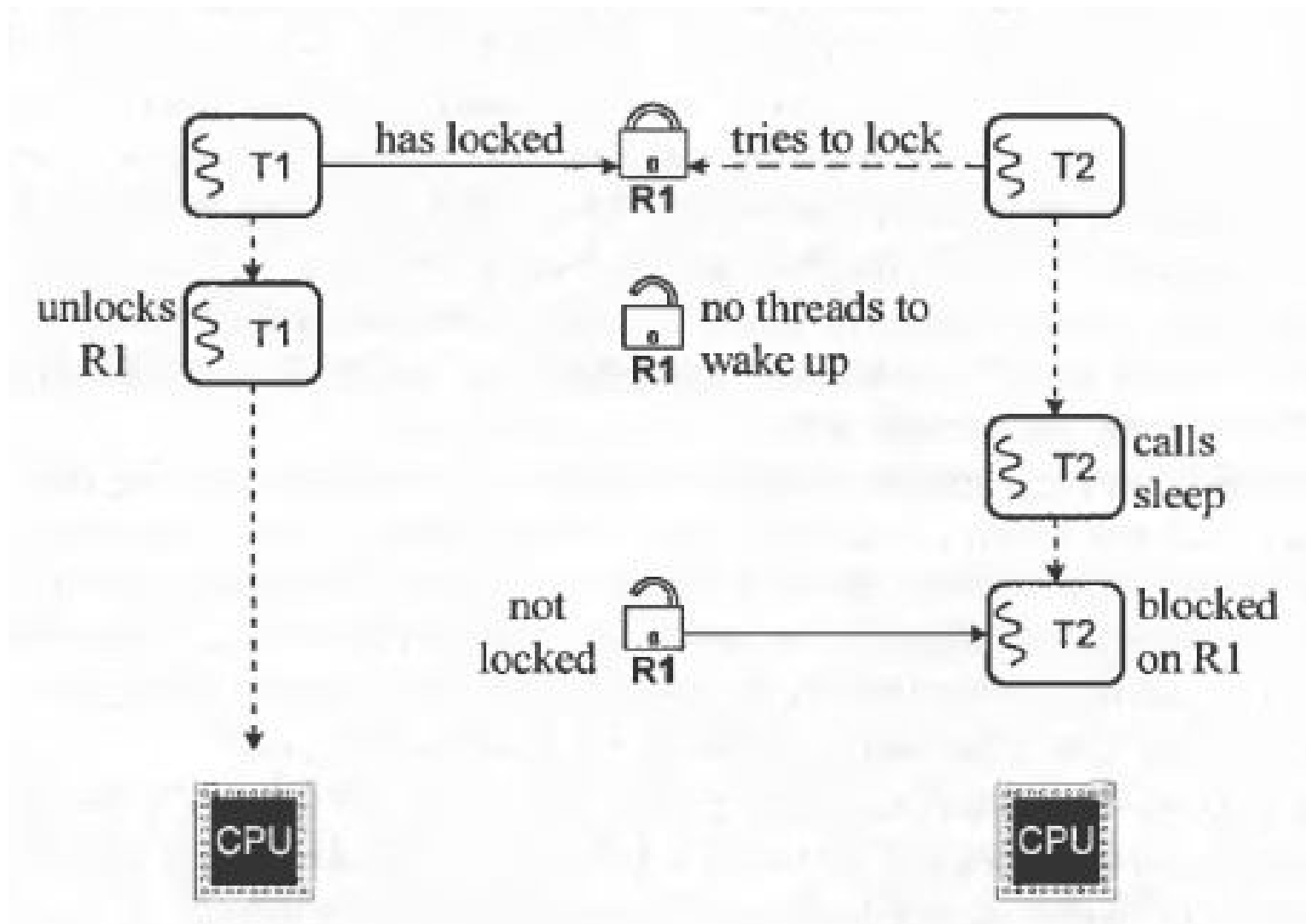


Figure 7-6. The lost wakeup problem.

Thundering herd problem

- **Thundering Herd problem**

- On a multiprocessor, if several threads were locked the resource
- Waking them all may cause them to be simultaneously scheduled on different processors
- and they would all fight for the same resource again.

- **Starvation**

- Even if only one thread was blocked on the resource, there is still a time delay between its waking up and actually running.
- In this interval, an unrelated thread may grab the resource causing the awakened thread to block again. If this happens frequently, it could lead to starvation of this thread.
- This problem is not as acute on a uniprocessor, since by the time a thread runs, whoever had locked the resource is likely to have released it.

Case Studies

Linux Synchronization

- **Prior to kernel Version 2.6, disables interrupts to implement short critical sections**
- **Version 2.6 and later, fully preemptive**
- **Linux provides:**
 - **semaphores**
 - **spinlocks**
 - **reader-writer versions of both**
 - **Atomic integers**
- **On single-cpu system, spinlocks replaced by enabling and disabling kernel preemption**

Linux Synchronization

- **Atomic variables**

`atomic_t` is the type for atomic integer

- **Consider the variables**

`atomic_t counter;`

`int value;`

<i>Atomic Operation</i>	<i>Effect</i>
<code>atomic_set(&counter,5);</code>	<code>counter = 5</code>
<code>atomic_add(10,&counter);</code>	<code>counter = counter + 10</code>
<code>atomic_sub(4,&counter);</code>	<code>counter = counter - 4</code>
<code>atomic_inc(&counter);</code>	<code>counter = counter + 1</code>
<code>value = atomic_read(&counter);</code>	<code>value = 12</code>

Pthreads synchronization

- **Pthreads API is OS-independent**
- **It provides:**
 - mutex locks
 - condition variables
- **Non-portable extensions include:**
 - read-write locks
 - spinlocks

Synchronization issues in xv6 kernel

Difference approaches

- **Pros and Cons of locks**
 - Locks ensure serialization
 - Locks consume time !
- **Solution – 1**
 - One big kernel lock
 - Too enefficient
- **Solution – 2**
 - One lock per variable
 - Often un-necessary, many data structures get manipulated in once place, one lock for all of them may work
- **Problem: ptable.lock for the entire array and every element within**
 - Alternatively: one lock for array, one lock per array entry

Three types of code

- **System calls code**
 - Can it be interruptible?
 - If yes, when?
- **Interrupt handler code**
 - Disable interrupts during interrupt handling or not?
 - Deadlock with iderw ! - already seen
- **Process's user code**
 - Ignore. Not concerned with it now.

Interrupts enabling/disabling in xv6

- **Holding every spinlock disables interrupts!**
- **System call code or Interrupt handler code won't be interrupted if**
 - **The code path followed took at least once spinlock !**
 - **Interrupts disabled only on that processor!**
- **Acquire calls `pushcli()` before `xchg()`**
- **Release calls `popcli()` after `xchg()`**

Memory ordering

- **Compiler may generate machine code for out-of-order execution !**
- **Processor pipelines can also do the same!**
- **This often improves performance**
- **Compiler may reorder 4 after 6 --> Trouble!**
- **Solution: Memory barrier**
 - **__sync_synchronize(), provided by GCC**
 - **Do not reorder across this line**
 - **Done only on acquire and release()**

- **Consider this**

```
1)l = malloc(sizeof *l);  
2)l->data = data;  
3)acquire(&listlock);  
4)l->next = list;  
5)list = l;  
6)release(&listlock);
```

Lost Wakeup?

- **Do we have this problem in xv6?**
- **Let's analyze again!**
 - The race in `acquiresleep()`'s call to `sleep()` and `releasesleep()`
- **T1 holding lock, T2 willing to acquire lock**
 - Both running on different processor
 - Or both running on same processor
 - What happens in both scenarios?
- **Introduce a T3 and T4 on each of two different processors. Now how does the scenario change?**
- **See page 69 in xv6 book revision-11.**

Code of sleep()

```
if(lk != &ptable.lock){  
    acquire(&ptable.lock);  
    release(lk);  
}
```

- **Why this check?**
- **Deadlock otherwise!**
- **Check: wait() calls with ptable.lock held!**

Exercise question : 1

Sleep has to check `lk != &ptable.lock` to avoid a deadlock

Suppose the special case were eliminated by replacing

```
if(lk != &ptable.lock){  
    acquire(&ptable.lock);  
    release(lk);  
}
```

with

```
release(lk);  
acquire(&ptable.lock);
```

Doing this would break sleep. How?

,

bget() problem

- **bget() panics if no free buffers!**
- **Quite bad**
- **Should sleep !**
- **But that will introduce many deadlock problems. Which ones ?**

iget() and ilock()

- **iget()** does not hold lock on inode
- **ilock()** does
- **Why this separation?**
 - Performance? If you want only “read” the inode, then why lock it?
- **What if iget() returned the inode locked?**

Interesting cases in namex()

```
while((path = skipelem(path, name)) != 0){
    ilock(ip);
    if(ip->type != T_DIR){
        iunlockput(ip);
        return 0;
    }
    if(nameiparent && *path == '\0'){
        // Stop one level early.
        iunlock(ip);
        return ip;
    }
}
```

```
if((next = dirlookup(ip, name, 0)) == 0){
    iunlockput(ip);
    return 0;
}
iunlockput(ip);
ip
}
```

--> only after obtaining next from dirlookup() and iget() is the lock released on ip;

-> lock on next obtained only after releasing the lock on ip. Deadlock possible if next was “.”

Xv6

Interesting case of holding and releasing
ptable.lock in scheduling

One process acquires, another releases!

Giving up CPU

- **A process that wants to give up the CPU**
 - must acquire the process table lock `ptable.lock`
 - release any other locks it is holding
 - update its own state (`proc->state`),
 - and then call `sched()`
- **Yield follows this convention, as do sleep and exit**
- **Lock held by one process P1, will be released another process P2 that starts running after `sched()`**
 - remember P2 returns either in `yield()` or `sleep()`
 - In both, the first thing done is releasing `ptable.lock`

Interesting race if ptable.lock is not held

- Suppose P1 calls yield()
- Suppose yield() does not take ptable.lock
 - Remember yield() is for a process to give up CPU
- Yield sets process state of P1 to RUNNABLE
- Before yield's sched() calls switch()
- Another processor runs scheduler() and runs P1 on that processor
- Now we have P1 running on both processors!
- P1 in yield taking ptable.lock prevents this

Homework

- **Read the version-11 textbook of xv6**
- **Solve the exercises!**

Scheduling

Abhijit A.M.
abhijit.comp@coep.ac.in

Credits: Slides from os-book.com

Necessity of scheduling

Multiprogramming

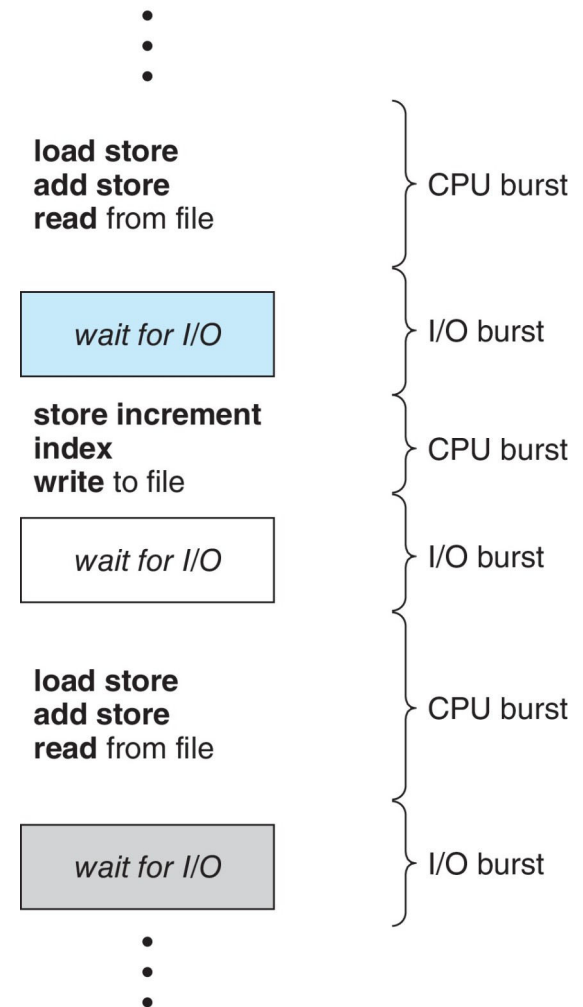
- Ability to run multiple programs at a time
- **Increase use of CPU**
 - CPU utilisation

CPU Scheduling

- **The task of selecting 'next' process/thread to execute on CPU and doing a context switch**
- **Scheduling algorithm**
 - **Criteria for selecting the 'next' process/thread and it's implementation**
- **Why is it important?**
 - **Affects performance !**
 - **Affects end user experience !**
 - **Involves money!**

Observation: CPU, I/O Bursts

- Process can 'wait' for an event (disk I/O, read from keyboard, etc.)
- During this period another process can be scheduled
- **CPU-I/O Burst Cycle:**
 - Process execution consists of a **cycle** of CPU execution and I/O wait
 - **CPU burst** followed by **I/O burst**
 - CPU burst distribution is of main concern



Let's understand the problem

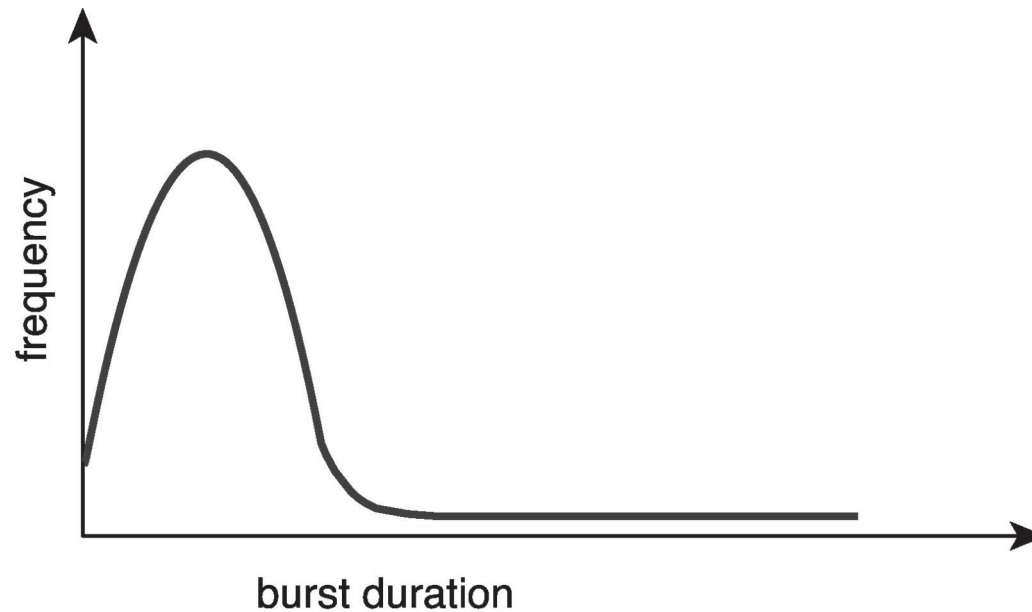
- **Programs have alternate CPU and I/O bursts**

- Some are CPU intensive
- Some are I/O intensive
- Some are mix of both

- **A C program example:**

```
f(int i, int j, int k) {  
    j = k * i; // CPU burst  
    scanf("%d", &i); // I/O burst  
    k = i * j; // CPU burst  
    printf("%d\n", k); // I/O burst  
    return k;  
}
```

CPU bursts: observation



Large number of short bursts

Small number of longer
bursts

Scheduler, what does it do?

- **From a list of processes, ready to run**
 - **Selects a process for execution**
 - **Allocates a CPU to the process for execution**
 - **Does “context switch”**
 - **Context: Set of registers**
 - **Switch from context of one process to another process**
 - **May be done like this: P1 context -> scheduler context -> P2 context**

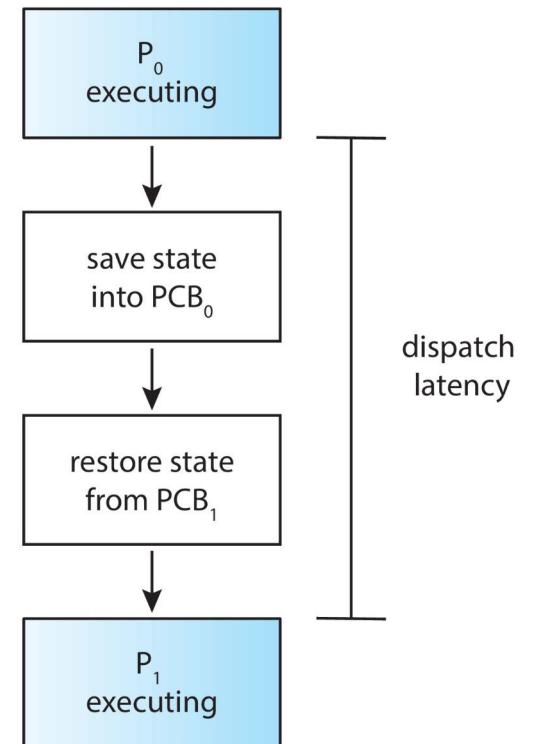
When is scheduler invoked?

- 1) **Process Switches from running to waiting state**
 - Waiting for I/O, etc.
- 2) **Switches from running to ready state**
 - E.g. on a timer interrupt
- 3) **Switches from waiting to ready**
 - I/O wait completed, now ready to run
- 4) **Terminates**

- **Scheduling under 1 and 4 is nonpreemptive**
- **All other scheduling is preemptive**
 - Consider access to shared data
 - Consider preemption while in kernel mode
 - Consider interrupts occurring during crucial OS activities

Dispatcher: A part of scheduler

- Gives control of the CPU to the process selected by the short-term scheduler; this involves:
 - switching context
 - switching to user mode
 - jumping to the proper location in the user program to restart that program
- **Dispatch latency**
 - time taken to stop one process and start another running
- **Xv6: `swtch()`, some tail end parts of `sched()`, `trap()`, `trapret()`**



Dispatcher in action on Linux

- **Run `vmstat 1 3`**
 - Means run `vmstat` 3 times at 1 second delay
- **In output, look at `CPU:cs`**
 - Context switches every second
- **Also for a process with pid 3323**
- **Run `cat /proc/3323/status`**
 - See
 - `voluntary_ctxt_switches`
--> Process left CPU
 - `nonvoluntary_ctxt_switches`
--> Process was preempted

Scheduling criteria

- **CPU utilization: Maximise**
 - keep the CPU as busy as possible. Linux: idle task is scheduled when no process to be scheduled.
- **Throughput : Maximise**
 - # of processes that complete their execution per time unit
- **Turnaround time : Minimise**
 - amount of time to execute a particular process
- **Waiting time : Minimise**
 - amount of time a process has been waiting in the ready queue
- **Response time : Minimise**
 - amount of time it takes from when a request was submitted until the first response is produced, not output (for time-sharing environment)

Calculations of different criteria

- **If you want to evaluate an algorithm practically, you need a proper workload !**
 - Processes with CPU and I/O bursts
 - Different durations of CPU bursts
 - Different durations of I/O bursts
 - How to do this programmatically?
 - How to ensure that after 2 seconds an I/O takes place?
 - Need periods when system will be “idle” – no process schedulable !

Calculations of different criteria

- **CPU Utilization**

- % time spent in doing 'useful' work
- What is useful work?
 - **On linux**
 - there is an "idle" thread, scheduled when no other task is RUNNABLE
 - Not running idle thread is productive work
 - Includes process + scheduling time + interrupts
 - **On other systems?**
 - Need to define
 - **On xv6**
 - We can say that time spent in the loop selecting a process is idle work

Calculations of different criteria

- **Throughput**

- # processes that complete execution per unit time
- Formula: total # processes completed / total time
- Simply divide by your total workload that completed by the time taken
- Depends on the workload as well. 'long' or 'short' processes.
- If too many short processes , then throughput may appear to be high, like 10s of processes per second
-

Calculations of different criteria

- **Turnaround time**
 - Amount of time required for one process to complete
 - For every process, note down the starting and ending time, difference is TA-time
 - For process P1 -> Time when process ended – time when process started
 - = Sum of time spent in (ready queue + running + waiting for I/O)
 - Do the average TA-time

Calculations of different criteria

- **Waiting time**
 - amount of time a process has been waiting in the *ready* queue.
 - **To be minimised.**
 - Part of turn around time
 - CPU scheduling does not affect waiting time in I/O queues, it affects time in ready queue

Scheduling Criteria

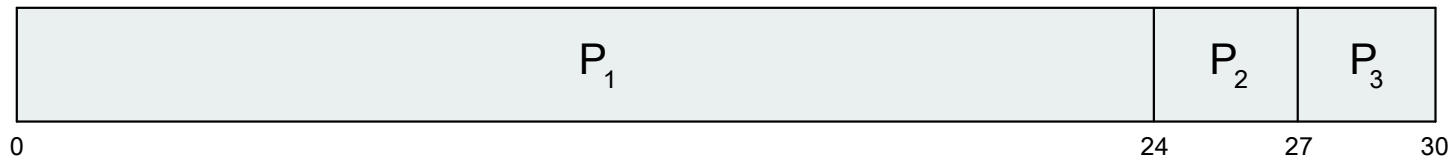
- **Response time**
 - amount of time it takes from when a request was submitted until the first response is produced, not output (for time-sharing environment).
 - **To be minimised.**
 - E.g. time between your press of a key , and that key being shown in screen

Challenges in implementing the scheduling algorithms

- **Not possible to know number of CPU and I/O bursts and the duration of each before the process runs !**
 - **Although when we do numerical problems around each algorithm, we assume some values for CPU and I/O bursts !**
-

GANTT chart

- A timeline chart showing the sequence in which processes get scheduled
- Used for analysing a scheduling algorithm



Scheduling Criteria: Differing requirements

- Different uses need different treatment of the scheduling criteria
- Knowing the workload is a challenge
- E.g. a desktop system
 - Response time is important
 - Minimize average response time Vs minimize variance in response time?
 - Sometimes fast, sometimes slow Vs Reasonable and predictable response time

Our discussion on scheduling algorithms

- **Analysis with only one CPU burst per process**
 - Ideally we should do for hundreds of CPU bursts
- **Only waiting time considered as criteria**

Scheduling Algorithms

First- Come, First-Served (FCFS) Scheduling

Process Burst Time

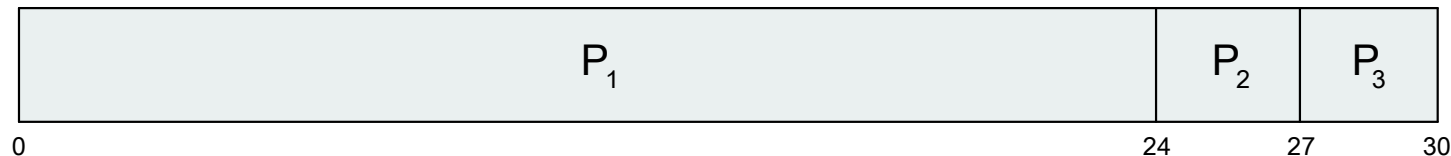
P1 24

P2 3

P3 3

Suppose that the processes arrive in the order: P1 , P2 , P3

The Gantt Chart for the schedule is:



Waiting time for P1 = 0; P2 = 24; P3 = 27

Average waiting time: $(0 + 24 + 27)/3 = 17$

Non Pre-emptive algorithm

FCFS Scheduling (Cont.)

Suppose that the processes arrive in the order:

P2 , P3 , P1

The Gantt chart for the schedule is:



Waiting time for P1 = 6; P2 = 0; P3 = 3

Average waiting time: $(6 + 0 + 3)/3 = 3$

Much better than previous case

FCFS: Convoy effect

- Consider one CPU-bound and many I/O-bound processes
- CPU bound process is scheduled, I/O bound processes are waiting in I/O queues
- I/O bound processes finish I/O and move to ready queue, and wait for CPU bound process to finish
 - I/O devices Idle
- CPU bound process over, goes for I/O. I/O bound processes run quickly, move to I/O queues again
 - CPU idle
- CPU bound process will run when it's ready to run
- Same process will repeat
- --> Lower CPU utilisation
 - Better if I/O bound processes run first

FCFS: further evaluation

- **Troublesome for interactive processes**
 - CPU bound process may hog CPU
 - Interactive process may not get a chance to run early and response time may be quite bad

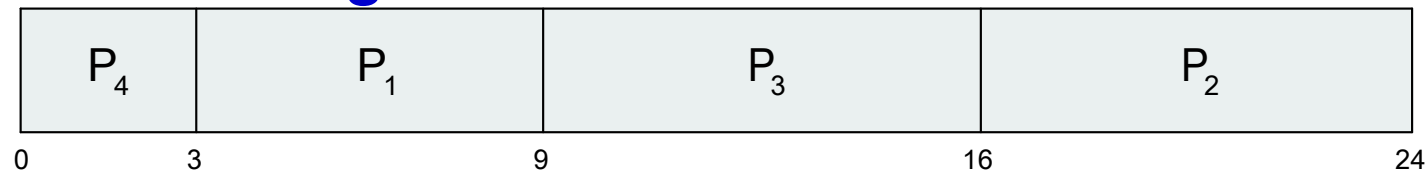
Shortest-Job-First (SJF) Scheduling

- Associate with each process the length of its next CPU burst
 - Use these lengths to schedule the process with the shortest time. Better name – **Shortest Next CPU Burst Scheduler**
- SJF is optimal – gives minimum average waiting time for a given set of processes
 - The difficulty is knowing the length of the next CPU request
 - Could ask the user – bad idea, unlikely to know!

Example of SJF

Process	Burst Time
P1	6
P2	8
P3	7
P4	3

SJF scheduling chart



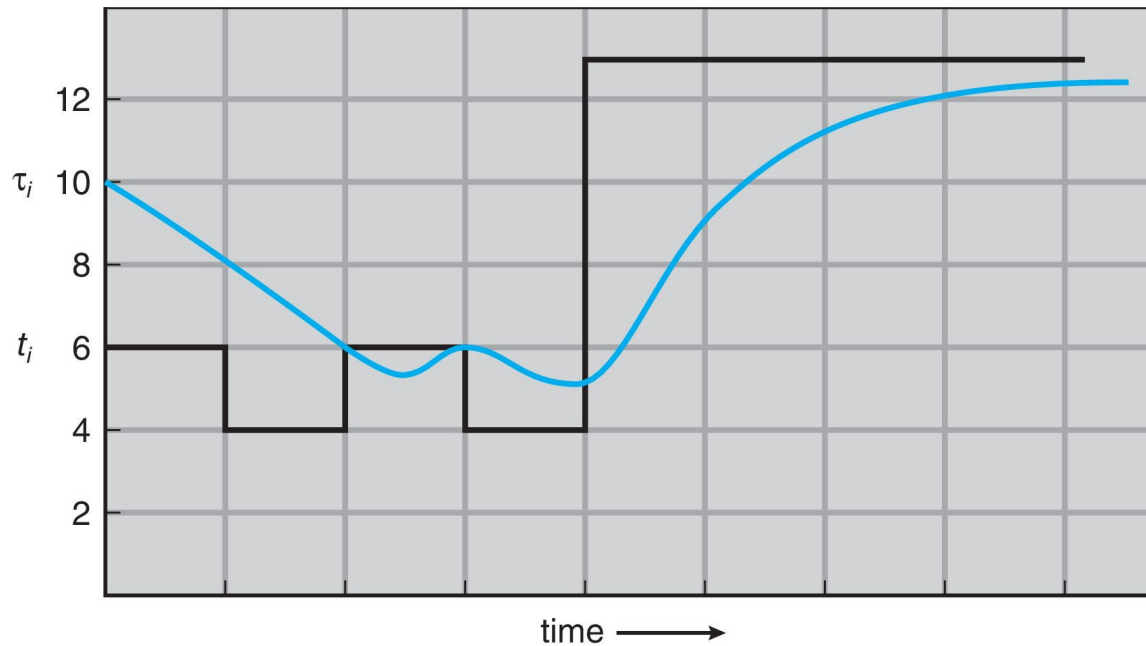
$$\text{Average waiting time} = (3 + 16 + 9 + 0) / 4 = 7$$

Determining Length of Next CPU Burst

- Not possible to implement SJF as can't know “next” CPU burst. Can only estimate the length – should be similar to the previous one
 - Then pick process with shortest predicted next CPU burst
- Can be done by using the length of previous CPU bursts, using exponential averaging
 1. t_n = actual length of n^{th} CPU burst
 2. τ_{n+1} = predicted value for the next CPU burst
 3. α , $0 \leq \alpha \leq 1$
 4. Define: $\tau_{n+1} = \alpha t_n + (1 - \alpha)\tau_n$.
- Commonly, α set to $\frac{1}{2}$
- Preemptive version called shortest-remaining-time-first

Prediction of the Length of the Next CPU Burst

- $A = 1/2$, $\tau_0 = 10$



CPU burst (t_i)	6	4	6	4	13	13	13	...	
"guess" (τ_i)	10	8	6	6	5	9	11	12	...

Examples of Exponential Averaging

- $\alpha = 0$

- $\tau_{n+1} = \tau_n$
- Recent history does not count

- $\alpha = 1$

- $\tau_{n+1} = \alpha t_n$
- Only the actual last CPU burst counts

- If we expand the formula, we get:

- $$\begin{aligned} \tau_{n+1} = & \alpha t_n + (1 - \alpha)\alpha t_{n-1} + \dots \\ & + (1 - \alpha)^j \alpha t_{n-j} + \dots \\ & + (1 - \alpha)^{n+1} \tau_0 \end{aligned}$$

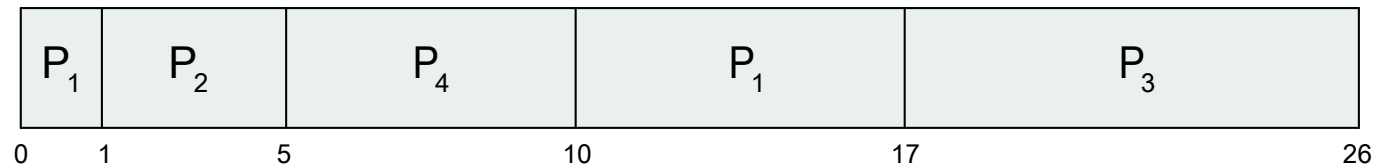
- Since both α and $(1 - \alpha)$ are less than or equal to 1, each successive term has less weight than its predecessor

Example of Shortest-remaining-time-first

Preemptive SJF = SRTF. Now we add the concepts of varying arrival times and preemption to the analysis

Process	Arrival Time	Burst Time
P1	0	8
P2	1	4
P3	2	9
P4	3	5

Preemptive SJF Gantt Chart



Average waiting time = $[(10-1)+(1-1)+(17-2)+(5-3)]/4 = 26/4 = 6.5$ msec

Round Robin (RR) Scheduling

- **Each process gets a small unit of CPU time (time quantum q), usually 10-100 milliseconds.**
 - After this time has elapsed, the process is preempted and added to the end of the ready queue.
- **If there are n processes in the ready queue and the time quantum is q , then each process gets $1/n$ of the CPU time in chunks of at most q time units at once.**
 - No process waits more than $(n-1)q$ time units.

Round Robin (RR) Scheduling

- **Timer interrupts every quantum to schedule next process**
- **Performance**
 - q large \Rightarrow FIFO
 - q small $\Rightarrow q$ must be large with respect to context switch, otherwise overhead is too high

Example of RR with Time Quantum = 4

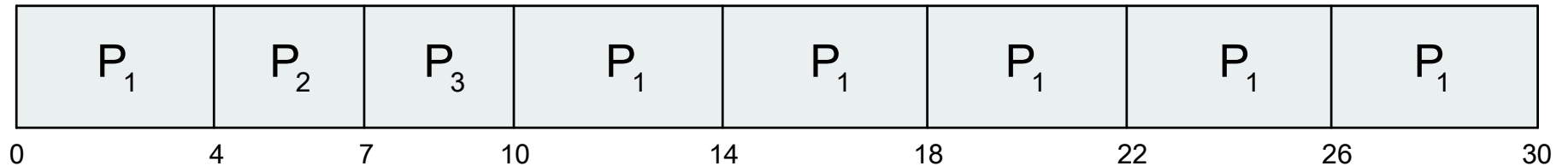
Process Burst Time

P_1 24

P_2 3

P_3 3

The Gantt chart is:

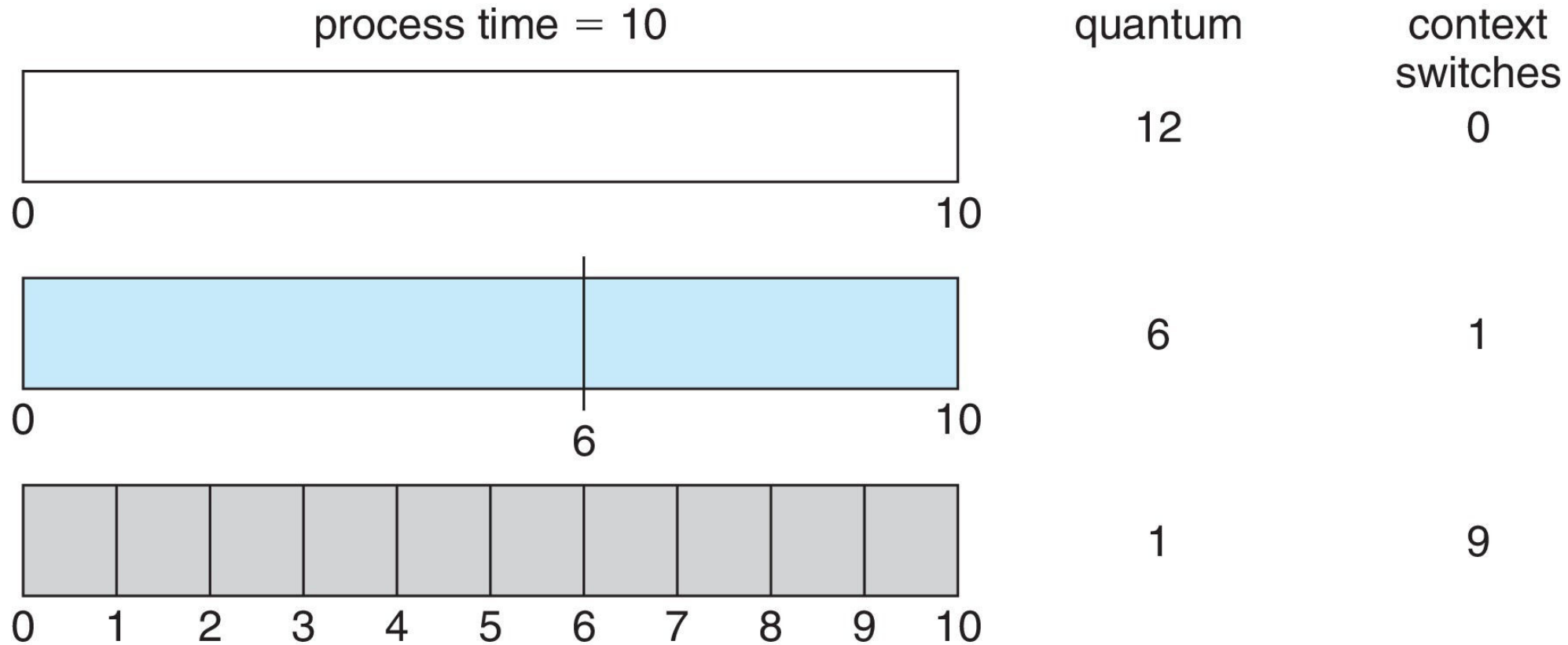


Typically, higher average turnaround than SJF, but better *response*

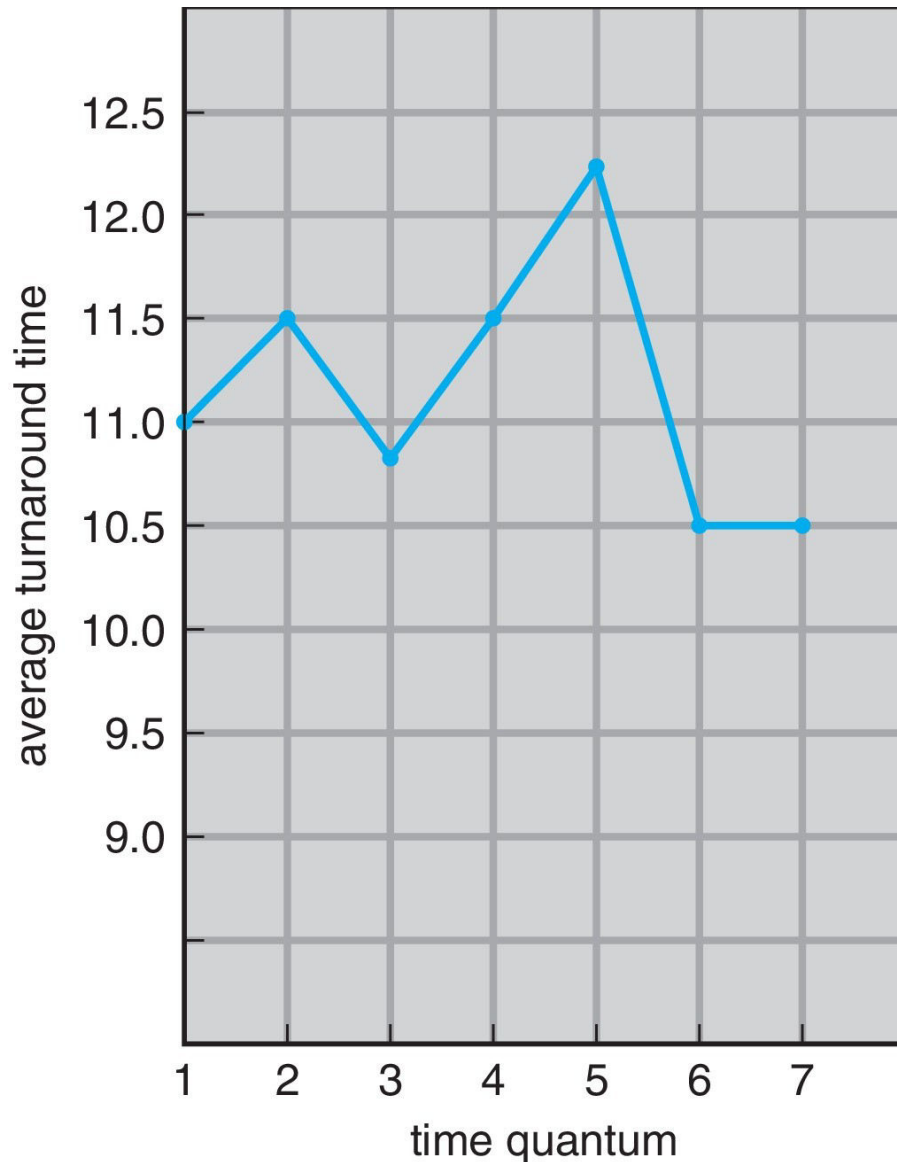
q should be large compared to context switch time

q usually 10ms to 100ms, context switch < 10 usec

Time Quantum and Context Switch Time



Turnaround Time Varies With The Time Quantum



process	time
P_1	6
P_2	3
P_3	1
P_4	7

80% of CPU bursts should be shorter than quantum

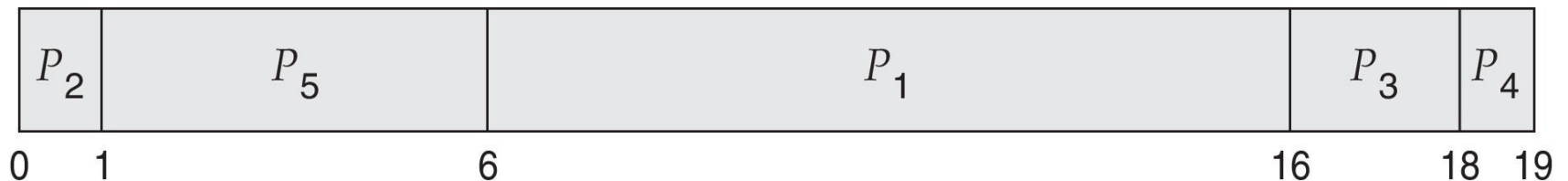
Priority Scheduling

- A priority number (integer) is associated with each process
- The CPU is allocated to the process with the highest priority (smallest integer \equiv highest priority)
 - Preemptive (timer interrupt, more time for more priority)
 - Nonpreemptive (no timer interrupt, just schedule process with highest priority)
- SJF is priority scheduling where priority is the inverse of predicted next CPU burst time
- Problem \equiv Starvation – low priority processes may never execute
- Solution \equiv Aging – as time progresses increase the priority of the process
-

Example of Priority Scheduling

<u>Process</u>	<u>Burst Time</u>	<u>Priority</u>
P_1	10	3
P_2	1	1
P_3	2	4
P_4	1	5
P_5	5	2

- Priority scheduling Gantt Chart



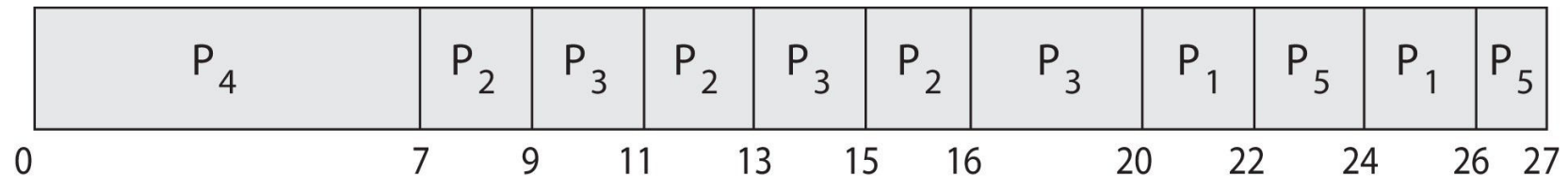
- Average waiting time = 8.2 msec

Priority Scheduling with Round-Robin

<u>Process</u>	<u>Burst Time</u>	<u>Priority</u>
P_1	4	3
P_2	5	2
P_3	8	2
P_4	7	1
P_5	3	3

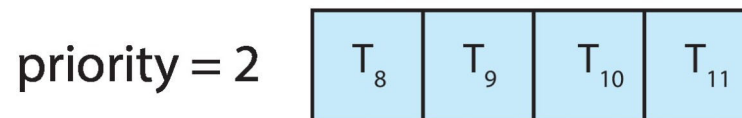
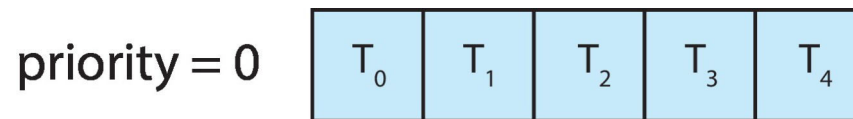
Run the process with the highest priority. Processes with the same priority run round-robin

Gantt Chart with 2 ms time quantum



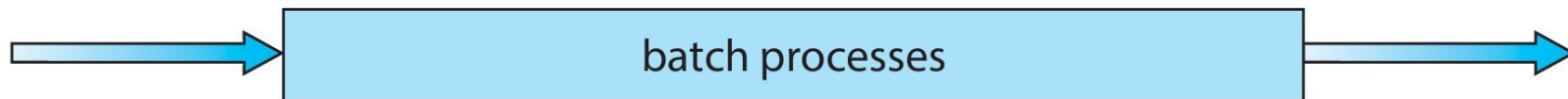
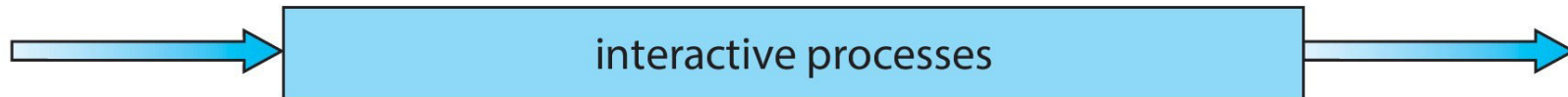
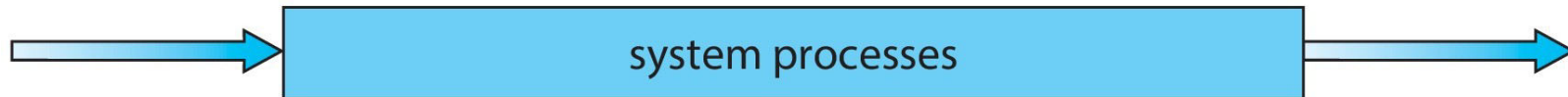
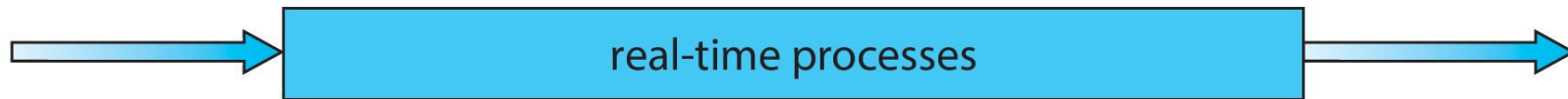
Multilevel Queue

- With priority scheduling, have separate queues for each priority.
- Schedule the process in the highest-priority queue!



Multilevel Queue

highest priority



lowest priority

Implementing multilevel queue

- **Processes need to have a priority**
 - Either modify `fork()/exec()` to have a priority
 - Or add a `nice()` system call to set priority
- **How to know the priority?**
 - The end user of the computer system needs to know this from needs of real life
 - E.g. on a database system, the database process will have a higher priority than other processes

Multilevel Feedback Queue

- A process can move between the various queues; aging can be implemented this way
- Multilevel-feedback-queue scheduler defined by the following parameters:
 - number of queues
 - scheduling algorithms for each queue
 - method used to determine when to upgrade a process
 - method used to determine when to demote a process
 - method used to determine which queue a process will enter when that process needs service

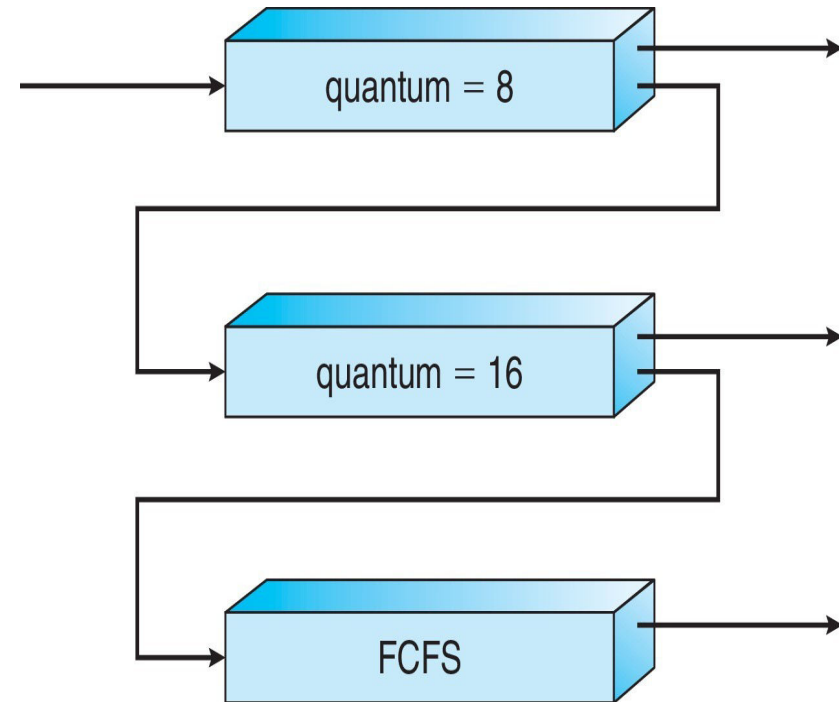
Example of Multilevel Feedback Queue

- Three queues:

- Q_0 – RR with time quantum 8 milliseconds
- Q_1 – RR time quantum 16 milliseconds
- Q_2 – FCFS

- Scheduling rules

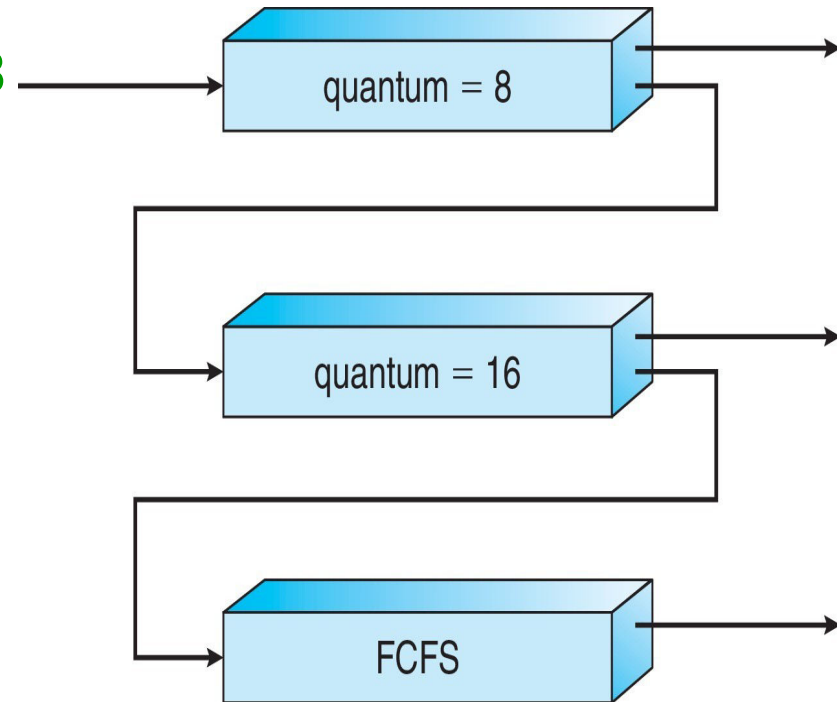
- Serve all processes in Q_0 first
- Only when Q_0 is empty, serve processes in Q_1
- Only when Q_0 and Q_1 are empty, serve processes in Q_2



Example of Multilevel Feedback Queue

- **Scheduling**

- A new job enters queue Q_0
 - When it gains CPU, job receives 8 milliseconds
 - If it does not finish in 8 milliseconds, job is moved to queue Q_1
- At Q_1 job receives 16 additional milliseconds
 - If it still does not complete, it is preempted and moved to queue Q_2
- To prevent starvation, move a process from lower-priority queue to higher priority queue after it has waited for too long



Thread Scheduling

- Distinction between user-level and kernel-level threads
- When threads supported, threads scheduled, not processes
- Many-to-one and many-to-many models, thread library schedules user-level threads to run on LWP
 - Known as **process-contention scope (PCS)** since scheduling competition is within the process
 - Typically done via priority set by programmer
- Kernel thread scheduled onto available CPU is **system-contention scope (SCS)** – competition among all threads in system

Pthread Scheduling

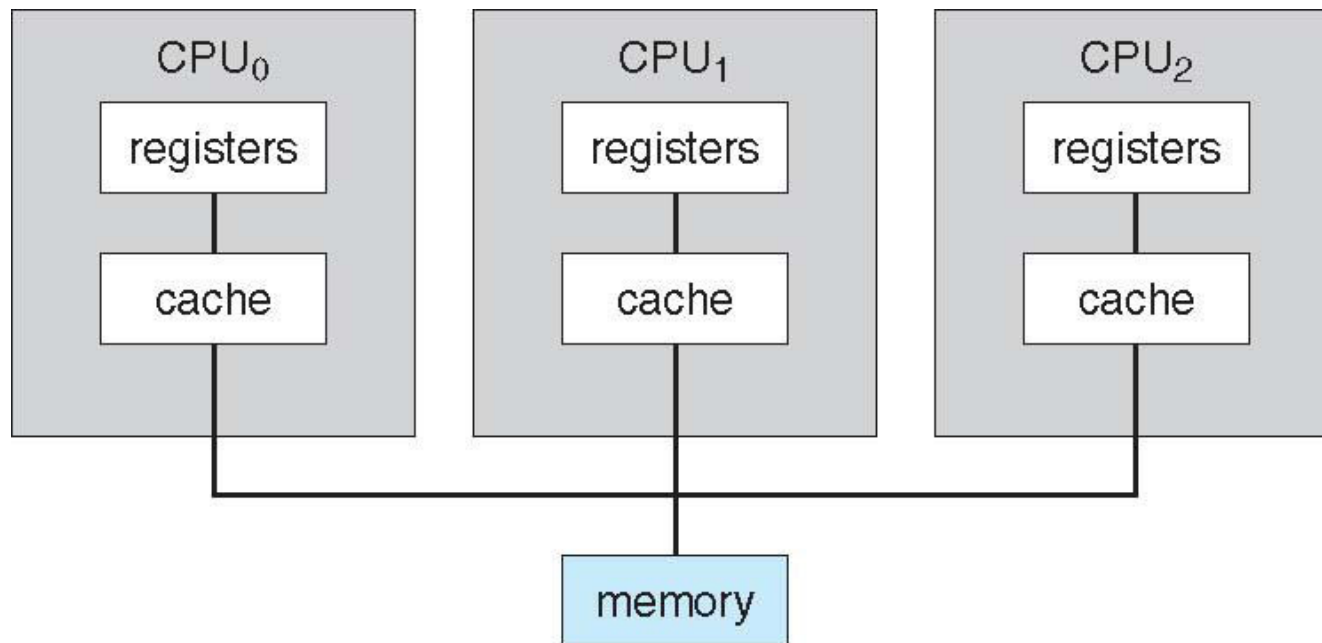
- **PTHREAD_SCOPE_PROCESS** schedules threads using PCS scheduling
- **PTHREAD_SCOPE_SYSTEM** schedules threads using SCS scheduling
- Linux and macOS only allow **PTHREAD_SCOPE_SYSTEM**
- Let's see a Demo using a program

Multi Processor Scheduling

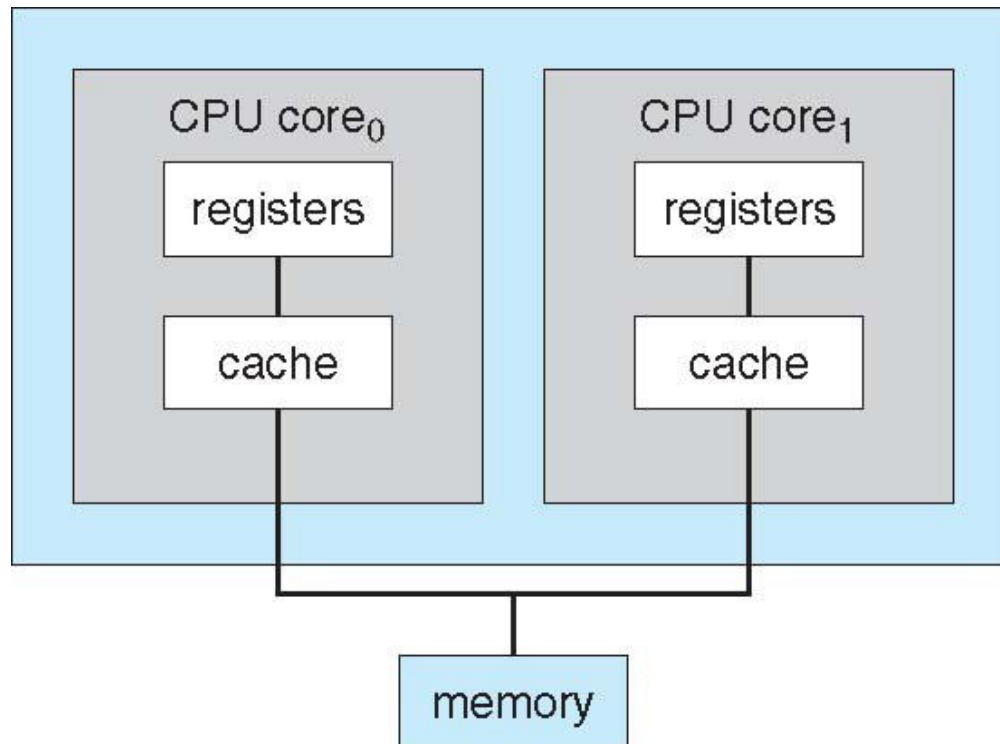
Multiprocessor systems

- **Each processor has separate set of registers**
 - All: eip, esp, cr3, eax, ebx, etc.
- **Each processor runs independently of others**
- **Main difference is in how do they access memory**

Symmetric multiprocessing (SMP)



Multicore systems (also SMP)



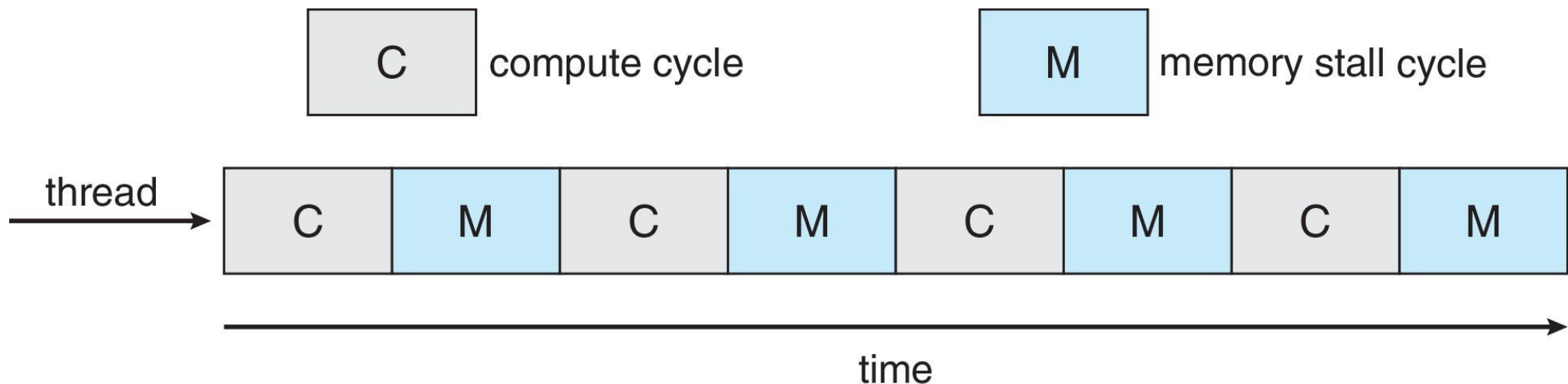
No difference from the perspective of OS. The hardware ensures that OS sees multiple processors and not multiple cores.

Booting multi-processor systems

- **X86 + xv6**
 - One processor starts, runs BIOS, loads kernel, initializes kernel data structures
 - **Mpinit()** : first processor scans certain memory addresses for information about other processors and obtains configuration and configures them
 - **Startothers()** : First processor initializes kernel data structures for other processors , and ensures that other processors start in scheduler() and starts them
 - First processor “starts” other processors

Multicore Processors

- Recent trend to place multiple processor cores on same physical chip
- Faster and consumes less power
- Multiple threads per core also growing
 - Takes advantage of memory stall to make progress on another thread while memory retrieve happens



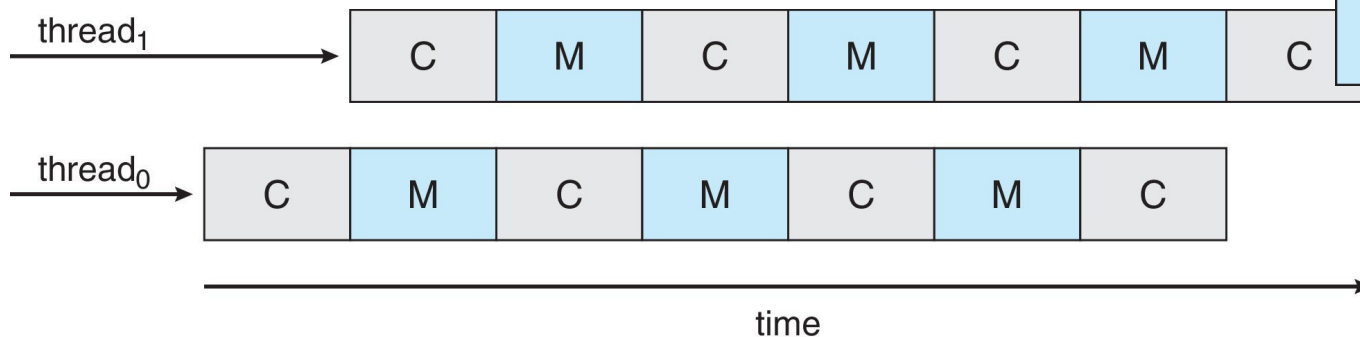
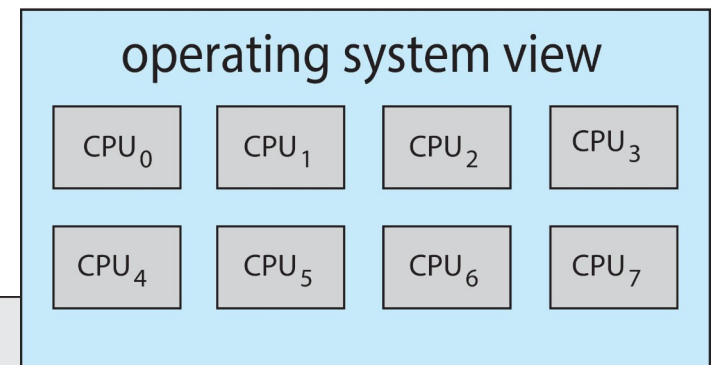
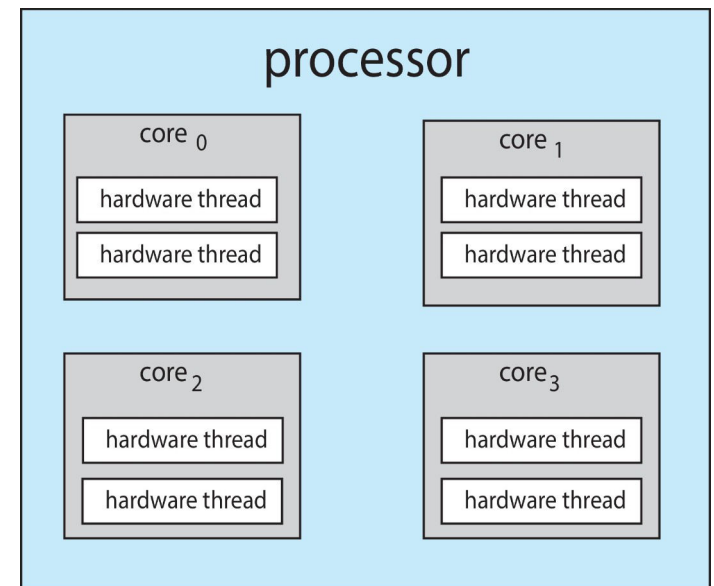
Multithreaded Multicore System

Each core has > 1 hardware threads.

Chip-multithreading (CMT) assigns each core multiple hardware threads. (Intel refers to this as **hyperthreading**.)

On a quad-core system with 2 hardware threads per core, the operating system sees 8 logical processors.

If one thread has a memory stall, switch to another thread!



Non-Uniform Memory Architecture (NUMA)

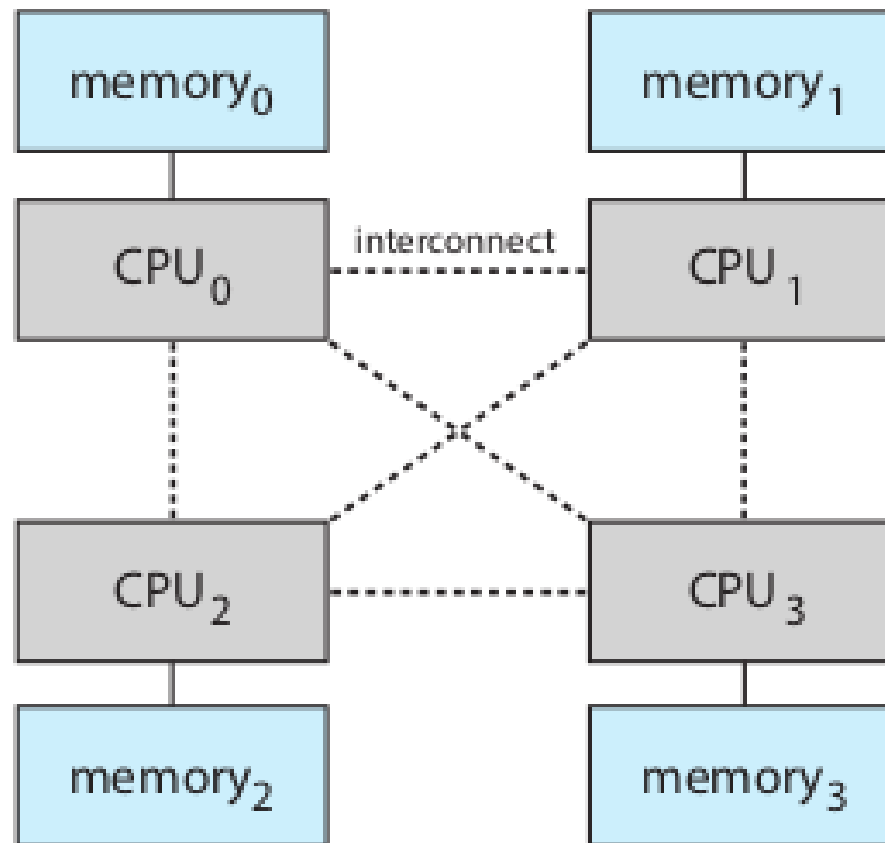


Figure 1.10 NUMA multiprocessing architecture.

More on SMP systems

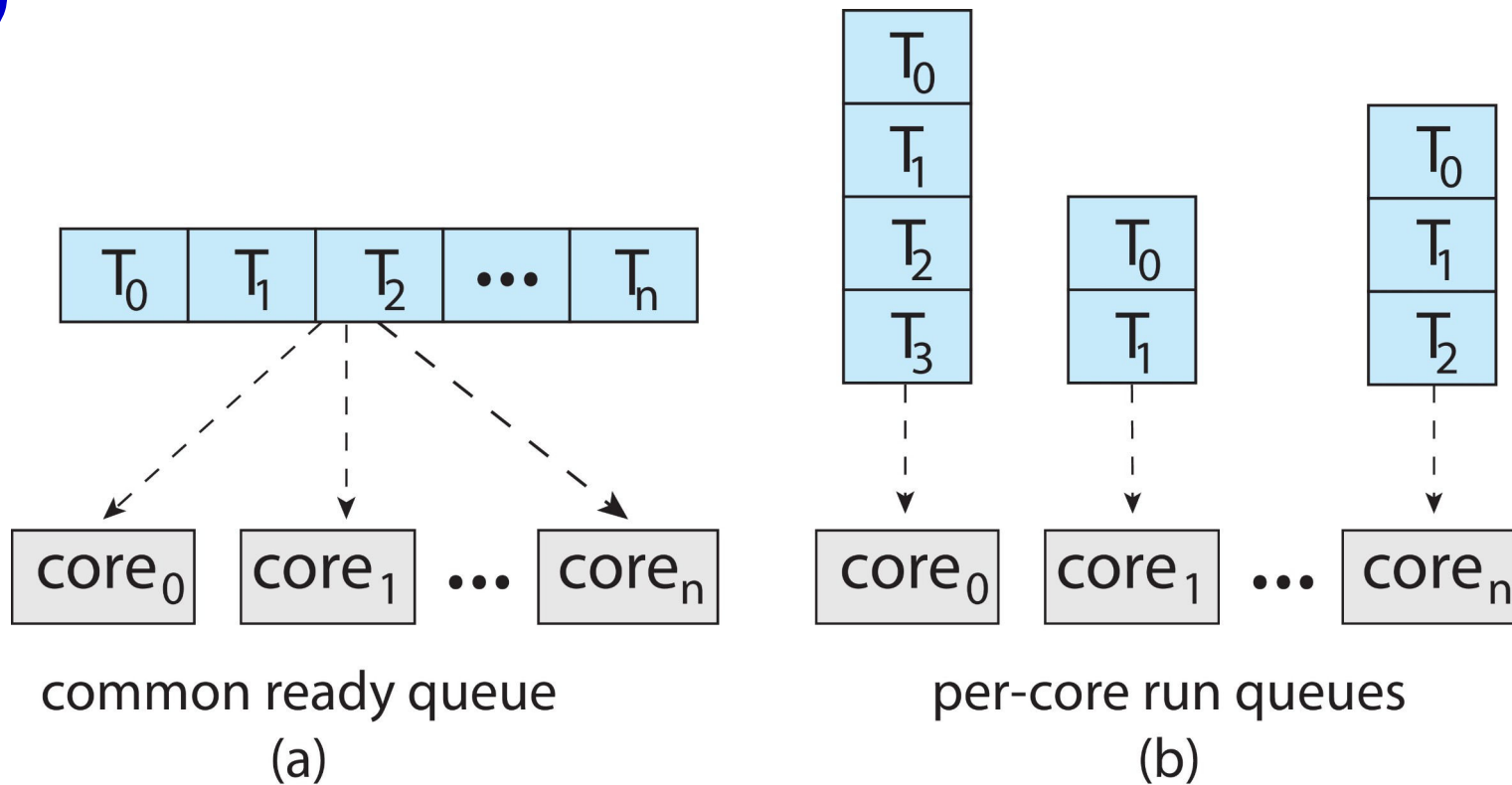
- **During booting – each CPU needs to be turned on**
 - Special I/O instructions writing to particular ports
 - See `lapicstartap()` in xv6
 - Need to setup CR3 on each processor
 - Segmentation, Page tables are shared (same memory for all CPUs)
- **All processors will keep running independently of each other**
- **Different interrupts on each processor – each needs IDT setup**
- **Each processor will be running processes, interrupt handlers, syscalls**
- **Synchronization problems !**
- **How to do scheduling ?**

Multiple-Processor Scheduling

- **CPU scheduling more complex when multiple CPUs are available**
- **Multiprocess may be any one of the following architectures:**
 - **Multicore CPUs**
 - **Multithreaded cores**
 - **NUMA systems**
 - **Heterogeneous multiprocessing**

Multiple-Processor Scheduling

- Symmetric multiprocessing (SMP) is where each processor is self scheduling.
- All threads may be in a common ready queue (a)
- Each processor may have its own private queue of threads (b)



Multiple-Processor Scheduling – Load Balancing

- If SMP, need to keep all CPUs loaded for efficiency
- Load balancing attempts to keep workload evenly distributed
- Push migration – periodic task checks load on each processor, and if found pushes task from overloaded CPU to other CPUs
- Pull migration – idle processors pulls waiting task from busy processor

Multiple-Processor Scheduling – Load Balancing

- When a thread has been running on one processor, the cache contents of that processor stores the memory accesses by that thread.
- We refer to this as a thread having affinity for a processor (i.e. “**processor affinity**”)
- Load balancing may affect processor affinity as a thread may be moved from one processor to another to balance loads, yet that thread loses the contents of what it had in the cache of the processor it was moved off of.

Multiple-Processor Scheduling – Load Balancing

- **Soft affinity** – the operating system attempts to keep a thread running on the same processor, but no guarantees.
- **Hard affinity** – allows a process to specify a set of processors it may run on.

SMP in xv6

- Only one process queue
- No load balancing, no affinity
- A process may run any CPU burst /alloted-time-quantum on any processor randomly
- See the code of:
 - Startothers(), mpenter(), mpmain()
- Different scheduler's kernel stack on each processor
 - Done in startothers()
- Each processor calls scheduler() from mpmain()

End