SYSC 5701

Operating System Methods for Real-Time Applications

Access Control: PCP Winter 2014

Resource-Sharing Dependencies

- A job cannot proceed (is blocked) because of resource-sharing synchronization
- Resource-sharing requires mutually exclusive access to the resource
- Can cause priority inversions

Feb 11/14

Resources

- serially reusable "units" of resource
 - eg. binary semaphore has one unit
 - counting semaphore has count units
- grant mutually exclusive right to access a unit
- once a unit is granted to a job, must not be reused by other jobs until released
- Recall management of mutual exclusion "unit" in monitors!

Feb 11/14

Access to Resources

- job requests resource(s)
 → job "locks" the resource(s)
- lock is managed by o/s (kernel)
- if resource(s) not available job is blocked
- eventually, job is granted the resource(s) and is unblocked
- when finished with resource(s)
 → job "unlocks" resource(s) for reuse

Feb 11/14

Access (more)

constructs to enable

related material from earlier in course:

- semaphores, IPC
- monitors
- critical sections

mutual exclusion -

parts of programs that require locking

programs to control locking

and unlocking of resources

the desired effect

Result: task can have <u>interdependencies</u> when accessing resources

Access Control Protocol

- resource conflict:
 - two jobs require same resource type
 - jobs must contend for the resource
- access control protocol: set of rules for
 - 1. granting resources
 - 2. scheduling jobs requesting resources

5

Feb 11/14

2

Priority Inversion

- a higher priority job is prevented from executing by a lower priority job
 - the priority relationship is inverted!



Unbounded Priority Inversion

• duration of priority inversion is not a function of the time for low priority job to execute the relevant critical section



Worst Case Job Response Time

- preemption time: delay due to higher priority job
- execution time: time to do job's work
- blocking time: time spent blocked
 - hopefully, blocking time is a simple function of delays while lower-priority jobs execute critical sections
 - if not, then difficult to compute (unbounded)

Avoiding Unbounded Priority Inversion

1. disable preemption

Feb 11/14

- 2. priority inheritance protocol
- 3. priority ceiling protocol

Disable All Preemption

- disable preemption during critical sections
- effectively elevate job in critical section to highest priority (cannot be preempted)
- priority elevation only needed when higherpriority jobs are requesting the relevant critical section – in other cases, the lower priority job should be preemptable by higher-priority jobs
- OK if critical sections are <u>very short</u> relative to shortest deadlines

11

Disable Preemption for Unbounded PI Example

10



Feb 11/14

Variation on Disable Preemption: **Priority Ceiling Emulation Priority Ceiling Emulation** in critical section, job runs at priority = priority ceiling for the resource **Priority Ceiling:** - i.e. no job that might request access to the resource is able to run! • the priority ceiling of resource R_i is the highest priority of all jobs that require access to R_i at job in critical section disables all jobs that might access critical section any time during their operation at end of critical section, job returns to original • denote $\Pi(R_i)$ priority • Q: do any jobs with priority higher than Π (R_i) jobs at priority higher than the ceiling are still access R: ? eligible to run Feb 11/14 Feb 11/14 13 14

Priority Ceiling Emulation for Unbounded PI Example



Priority Ceiling Emulation <u>vs.</u> Disable Preemption Example 2



Priority Ceiling Emulation <u>vs.</u> Disable Preemption Example 1



Basic Priority Inheritance Protocol

- while a job_{low} is holding any resource: raise its priority to the highest priority of any job requesting any resource held by job_{low}
- dynamic → raise at time higher priority job requests the resource
- when unlock a resource: assign job_{low} the higher of (1) its original priority, or (2) the highest priority of a job requesting a resource held by job_{low}

Feb 11/14



19

Feb 11/14



Chain Blocking Example



Basic Priority Inheritance for Unbounded PI Example



Is Blocking a Function of Time to Execute Critical Sections?

- suppose *m* critical sections accessed by task τ
- job of τ can be blocked directly, at most, m times
 - not counting indirect blocking!
- if *n* tasks at lower priority than τ
 - job of τ can be blocked, at most, at one critical section in each of the n tasks
 - blocking time is bounded,
 but ... may suffer from "chain blocking" – blocks each time attempt to access a critical section

22

24

Feb 11/14

Chain Blocking for Disable Preemption & Priority Inheritance?

- still a problem!
- previous example:





Allocation Rule (con't)

if R is free:

Feb 11/14

- i. if priority of J at t > Λ(t), allocate R to J
 J does not access any of the held resources!
- ii. else: if J is the job holding the resource(s) whose priority ceiling = $\Pi(t)$, allocate R to J
 - Not possible for different jobs to hold resources with same priority ceiling! (see i. and iii.)

iii. otherwise: request denied and J is blocked

Allocation Rule (paraphrasing)

- a job <u>cannot</u> acquire a resource unless its priority is higher than the ceilings of all other resources currently acquired by other jobs
- if priority higher than ceilings, then job will not request access to any of the other active resources (by the definition of a ceiling!)
- when request denied and job is blocked, higher priority jobs <u>might still be able to acquire</u> <u>resource</u>! (deny access ≠ FIFO blocking)

32

34

Priority-Inheritance Rule

- while a job_{low} is holding any resource: raise its priority to the highest priority of any job requesting any resource held by job_{low}
- dynamic: at time t when a job J becomes blocked, the job J_{low} which blocks J inherits the current priority of J
- J_{low} executes at inherited priority until t' when it releases every resource whose priority ceiling is greater or equal to the inherited priority
- at t': priority of J_{low} falls to the higher of (1) its original priority, or (2) the priority of the highest job requesting one of the resources still held by J_{low}

Feb 11/14

33

31

Lowering Priority Scenario

- Suppose job_{low} holds resource 1 with priority ceiling π₃ which is then requested by job with priority π₄ (rules?)
 → raise job_{low} to priority π₄
- Now job_{low} acquires resource 2 with priority ceiling π₁ which is then requested by job with priority π₂ (rules?)
 → raise job_{low} to priority π₂
- Now job_{low} releases resource 1
 → what should be priority of job_{low} now? (rules?)
- Now job_{low} releases resource 2
 → what should be priority of job_{low} now? (rules?)
- What if resources released in other order? (rules?)

Feb 11/14

Feb 11/14

Recall Previous Deadlock Example



Another Example

 $\begin{aligned} \tau_0 &: \ \mathsf{L}(\mathsf{S}_0) \ \mathsf{U}(\mathsf{S}_0) \ \mathsf{L}(\mathsf{S}_1) \ \mathsf{U}(\mathsf{S}_1) \\ \tau_1 &: \ \mathsf{L}(\mathsf{S}_2) \ \mathsf{U}(\mathsf{S}_2) \\ \tau_2 &: \ \mathsf{L}(\mathsf{S}_2) \ \mathsf{L}(\mathsf{S}_1) \ \mathsf{U}(\mathsf{S}_1) \ \mathsf{U}(\mathsf{S}_2) \\ \Pi(\mathsf{S}_0) &= \ \pi(\tau_0) \ \Pi(\mathsf{S}_1) &= \ \pi(\tau_0) \ \Pi(\mathsf{S}_2) &= \ \pi(\tau_1) \end{aligned}$





Points Seen in Example

- a job is blocked by a (possibly nested) critical section of at most one lower priority job
- ceiling blocking occurs a job is prevented from entering a critical section by ceiling of an active resource \rightarrow not because the requested resource was busy !

– e.g. t_2 : τ_0 is blocked even though S_0 is free

Properties of Basic Priority Ceiling Protocol

• no deadlock !

Feb 11/14

- job blocks in at most one critical section blocking is bounded
 - no chain blocking \rightarrow shorter blocking bound than Priority Inheritance Protocol
- once acquire first resource, all resources needed will be available when requested

Feb 11/14	39	

Implementation of Basic **Priority Ceiling Protocol**

- don't need "lock" queues (e.g. semaphore queue)
- maintain queue of tasks that are ready-to-run or blocked – maintain in priority order task at head is current task

Why is a single queue sufficient?

 need list of active resources – ordered by ceiling priority (includes task that locked the resource, and highest priority of any task blocked waiting for the resource)

- lock and unlock manipulate queue and list
- need analysis of critical section use establish priority ceilings prior to run-time

Example Using Client / Server (similar to Liu text)

40



Feb 11/14



Basic Priority Inheritance Same Example



Response Comparison

