

Real-Time Scheduling (Part 2)

(Working Draft)

Insup Lee
Department of Computer and Information Science
School of Engineering and Applied Science
University of Pennsylvania
www.cis.upenn.edu/~lee/



CIS 541, Spring 2010

PERIODIC SERVERS FOR APERIODIC TASKS

Mixed Periodic and Aperiodic Task Systems

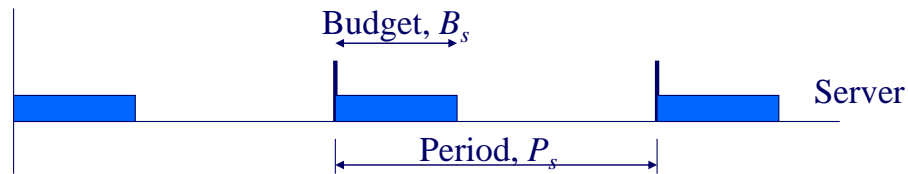
- Question: how to execute aperiodic tasks without violating schedulability guarantees given to periodic tasks?

Mixed Periodic and Aperiodic Task Systems

- Question: how to execute aperiodic tasks without violating schedulability guarantees given to periodic tasks?
- One Answer: Execute aperiodic tasks at lowest priority
 - Problem: Poor performance for aperiodic tasks

Mixed Periodic and Aperiodic Task Systems

- Idea: aperiodic tasks can be served by periodically invoked servers
- The server can be accounted for in periodic task schedulability analysis
- The server has a period P_s and a budget B_s
- Server can serve aperiodic tasks until budget expires
- Servers have different flavors depending on the details of when they are invoked, what priority they have, and how budgets are replenished



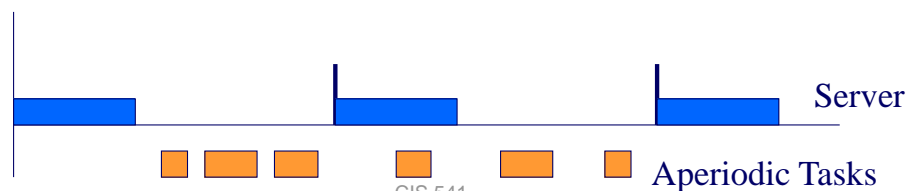
Spring '10

CIS 541

5

Mixed Periodic and Aperiodic Task Systems

- Idea: aperiodic tasks can be served by periodically invoked servers
- The server can be accounted for in periodic task schedulability analysis
- The server has a period P_s and a budget B_s
- Server can serve aperiodic tasks until budget expires
- Servers have different flavors depending on the details of when they are invoked, what priority they have, and how budgets are replenished



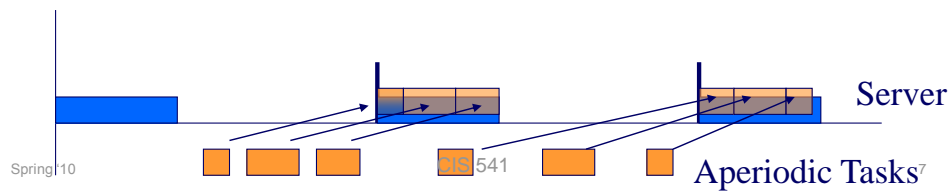
Spring '10

CIS 541

6

Mixed Periodic and Aperiodic Task Systems

- Idea: aperiodic tasks can be served by periodically invoked servers
- The server can be accounted for in periodic task schedulability analysis
- The server has a period P_s and a budget B_s
- Server can serve aperiodic tasks until budget expires
- Servers have different flavors depending on the details of when they are invoked, what priority they have, and how budgets are replenished

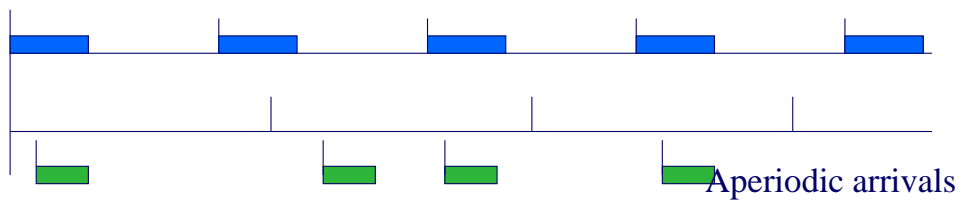


Polling Server

- Runs as a periodic task (priority set according to RM)
- Aperiodic arrivals are queued until the server task is invoked
- When the server is invoked it serves the queue until it is empty or until the budget expires then suspends itself
 - If the queue is empty when the server is invoked it suspends itself immediately.
- Server is treated as a regular periodic task in schedulability analysis

Example of a Polling Server

- Polling server:
 - Period $P_s = 5$
 - Budget $B_s = 2$
- Periodic task
 - $P = 4$
 - $C = 1.5$
- All aperiodic arrivals have $C=1$



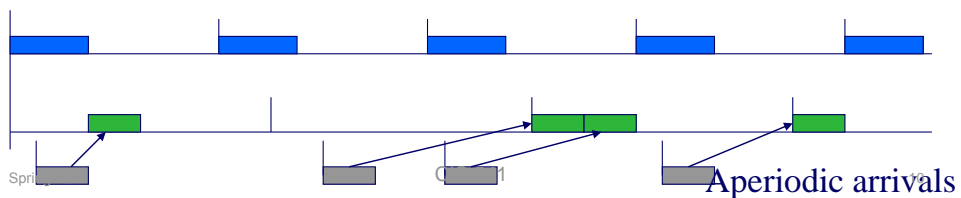
Spring '10

CIS 541

9

Example of a Polling Server

- Polling server:
 - Period $P_s = 5$
 - Budget $B_s = 2$
- Periodic task
 - $P = 4$
 - $C = 1.5$
- All aperiodic arrivals have $C=1$



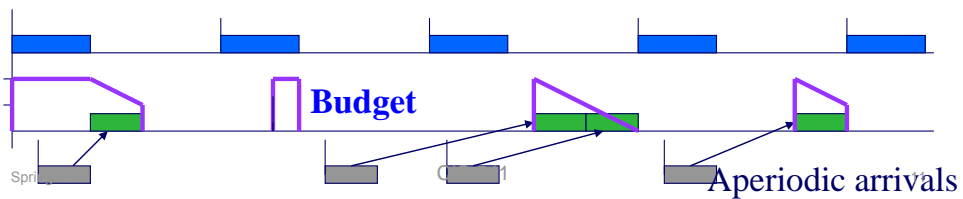
Spring

CIS 541

9

Example of a Polling Server

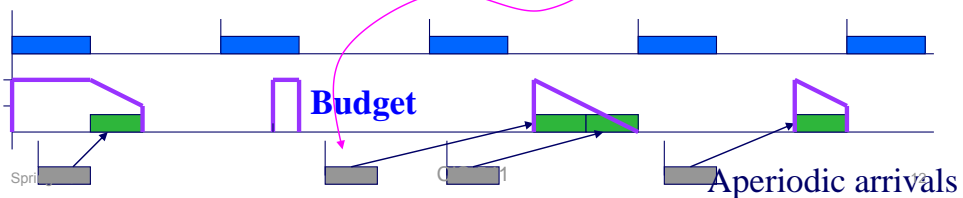
- Polling server:
 - Period $P_s = 5$
 - Budget $B_s = 2$
- Periodic task
 - $P = 4$
 - $C = 1.5$
- All aperiodic arrivals have $C=1$



Example of a Polling Server

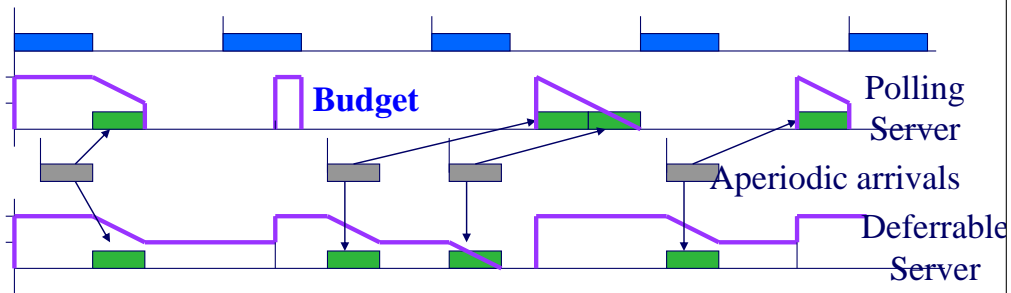
- Polling server:
 - Period $P_s = 5$
 - Budget $B_s = 2$
- Periodic task
 - $P = 4$
 - $C = 1.5$
- All aperiodic arrivals have $C=1$

Why not execute immediately?



Deferrable Server

- Keeps the balance of the budget until the end of the period
- Example (continued)

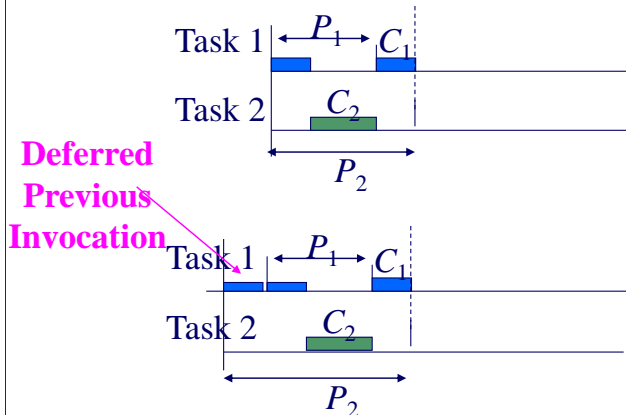


Spring '10

CIS 541

13

Worst-Case Scenario



$$U_p \leq \ln \left(\frac{2U_s + 1}{U_s + 1} \right)$$

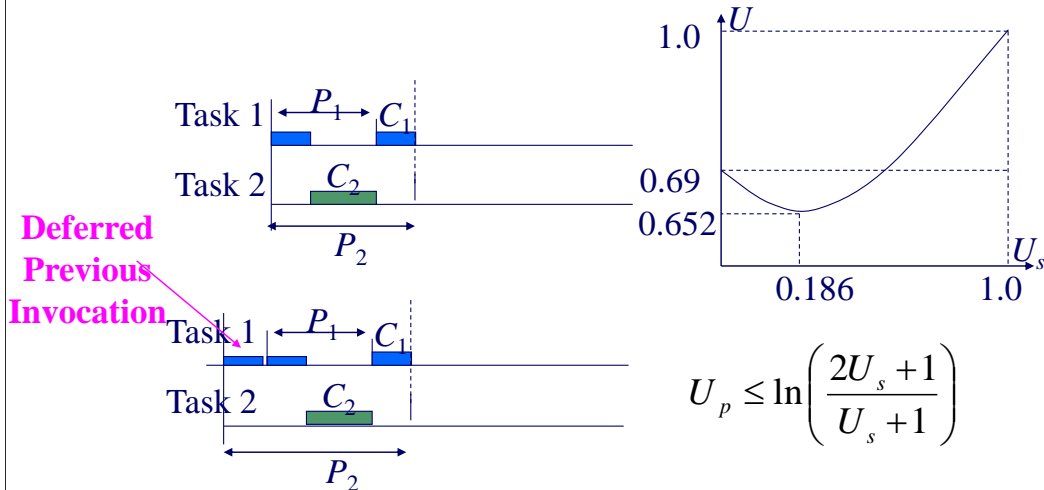
Exercise: Derive the utilization bound for a deferrable server plus one periodic task

Spring '10

CIS 541

14

Worst-Case Scenario



Exercise: Derive the utilization bound for a deferrable server plus one periodic task

Spring '10

CIS 541

15

Priority Exchange Server

- Like the deferrable server, it keeps the budget until the end of server period
- Unlike the deferrable server the priority slips over time: When not used, the priority is exchanged for that of the executing periodic task

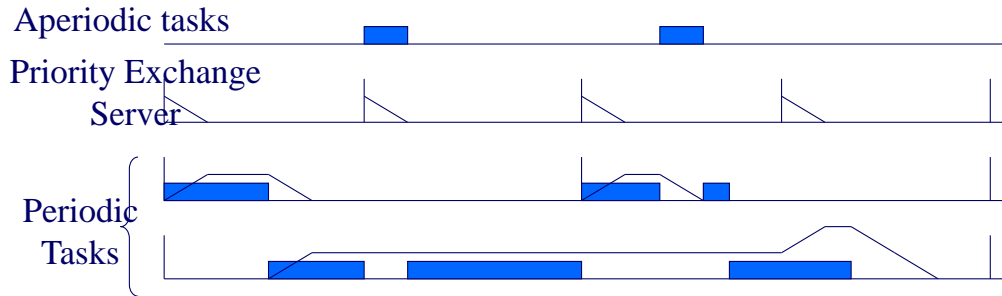
Spring '10

CIS 541

16

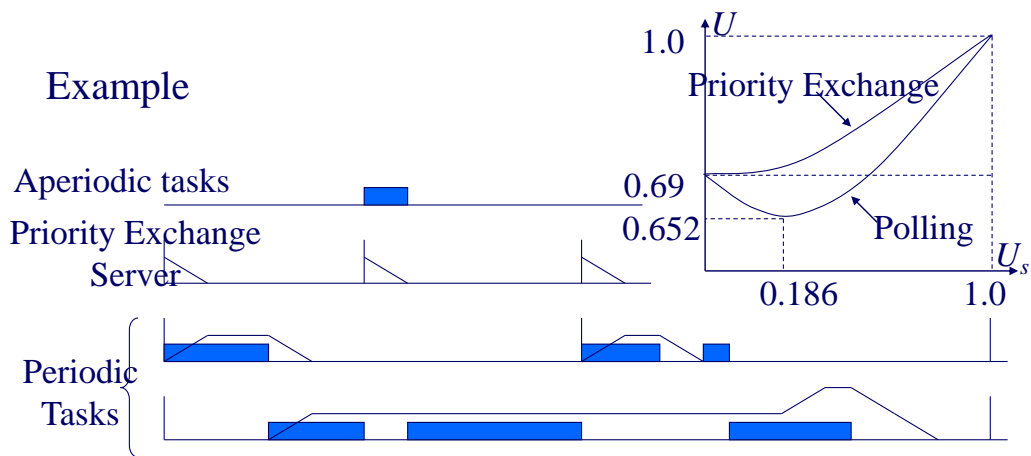
Priority Exchange Server

Example



Priority Exchange Server

Example



$$U_p \leq \ln\left(\frac{2}{U_s + 1}\right)$$

Sporadic Server

- Server is said to be *active* if it is in the *running* or *ready* queue, otherwise it is *idle*.
- When an aperiodic task comes and the budget is not zero, the server becomes active
- Every time the server becomes *active*, say at t_A , it sets replenishment time one period into the future, $t_A + P_s$ (but does not decide on replenishment amount).
- When the server becomes *idle*, say at t_I , set replenishment amount to capacity consumed in $[t_A, t_I]$

$$U_p \leq \ln\left(\frac{2}{U_s + 1}\right)$$

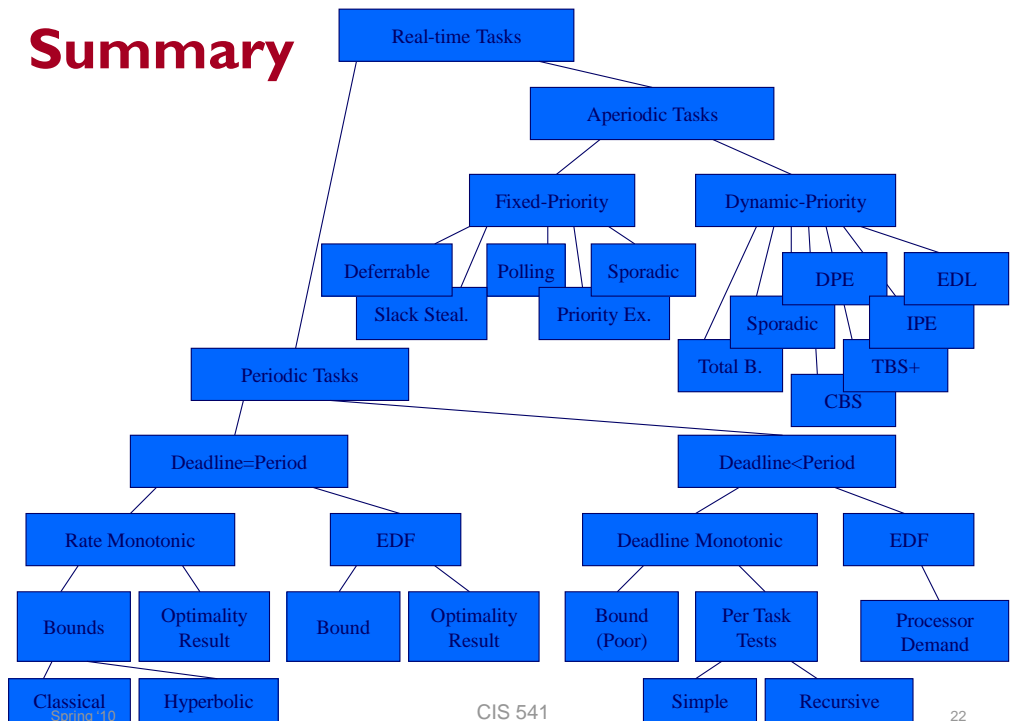
Slack Stealing Server

- Compute a slack function $A(t_s, t_f)$ that says how much total slack is available
- Admit aperiodic tasks while slack is not exceeded

Aperiodic Servers

Dynamic Priority

Summary



PRIORITY INVERSION

Resources and Blocking

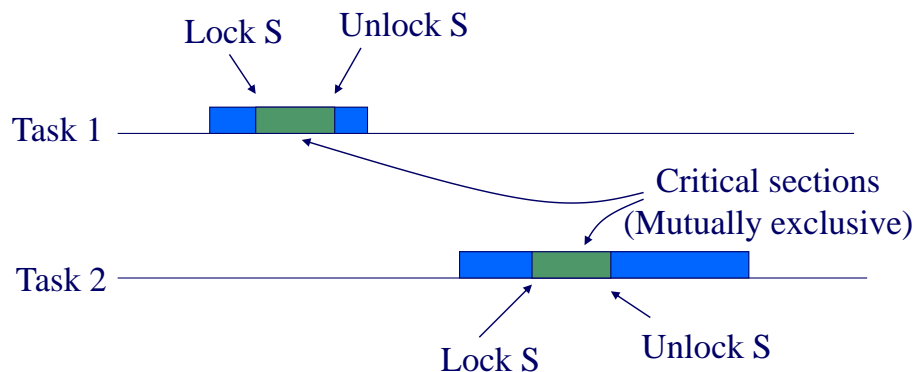
Priority Inheritance
Priority Ceiling
Slack Resource Policy

The Problem

- Tasks have synchronization constraints
 - Semaphores protect critical sections
- Blocking can cause a higher-priority task to wait on a **lower-priority** one to unlock a resource
 - Problem: In all previous derivations we assumed that a task can only wait for **higher-priority** tasks not **lower-priority** tasks
- Question
 - What is the maximum amount of time a higher-priority task can wait for a lower-priority task?
 - How to account for that time in schedulability analysis?

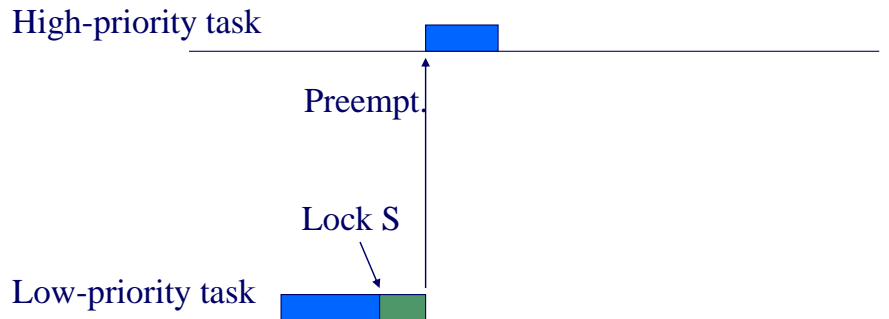
Mutual Exclusion Constraints

- Tasks that lock/unlock the same semaphore are said to have a mutual exclusion constraint



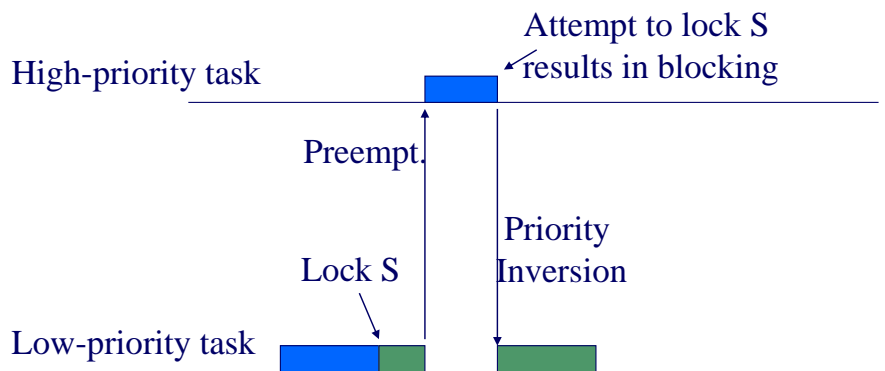
Priority Inversion

- Locks and priorities may be at odds. Locking results in priority inversion



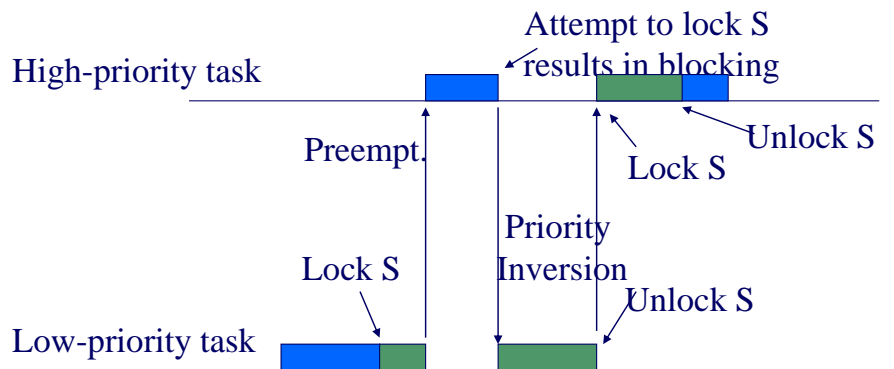
Priority Inversion

- Locks and priorities may be at odds. Locking results in priority inversion



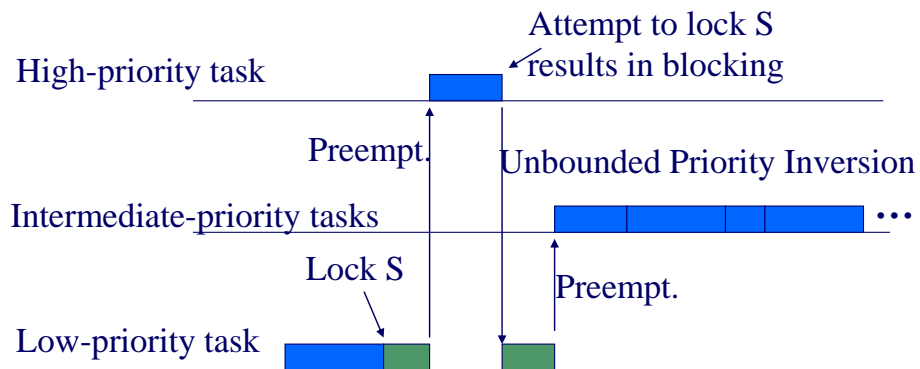
Priority Inversion

- How to account for priority inversion?



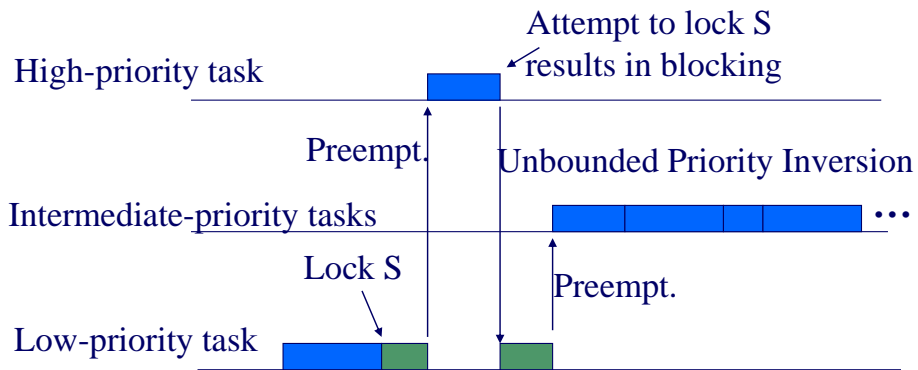
Unbounded Priority Inversion

- Consider the case below: a series of intermediate priority tasks is delaying a higher-priority one



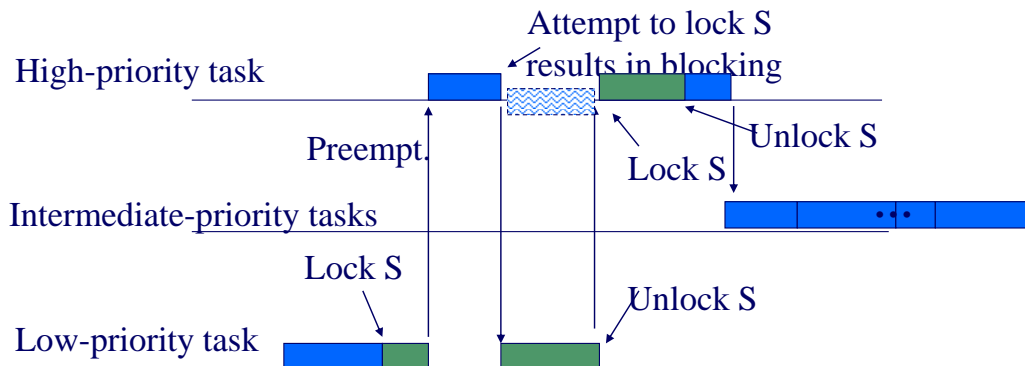
Unbounded Priority Inversion

- How to prevent unbounded priority inversion?



Priority Inheritance Protocol

- Let a task inherit the priority of any higher-priority task it is blocking

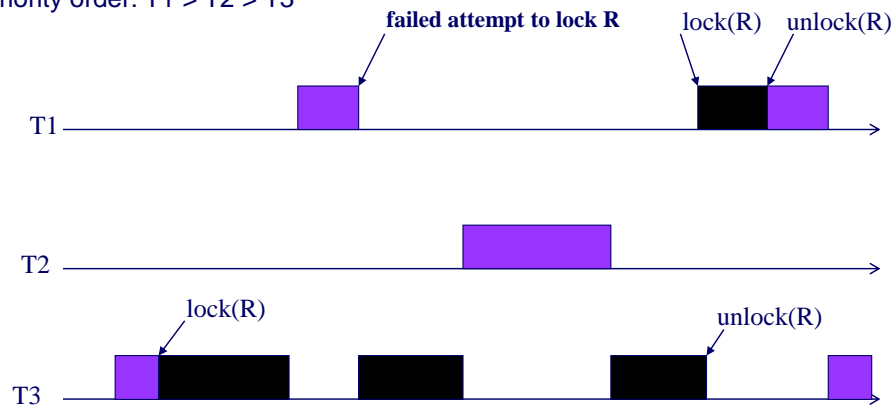


Priority Inversion and the MARS Pathfinder

- Landed on the Martian surface on July 4th, 1997
- Unconventional landing – bouncing into the Martian surface
- A few days later, not long after Pathfinder started gathering meteorological data, the spacecraft began experiencing total system reset, each resulting in losses of data
- What happened:
 - Pathfinder has an “information bus”
 - The meteorological data gathering task ran as an infrequent, low priority thread, and used the information bus to publish its data (while holding the mutex on bus).
 - A communication task that ran with medium priority.
 - It is possible for an interrupt to occur that caused (medium priority) communications task to be scheduled during the short interval of the (high priority) information bus thread was blocked waiting for the (low priority) meteorological data thread.
 - After some time passed, a watch dog timer goes off, noticing that the data bus has not been executed for some time, it concluded that something had gone really bad, and initiated a total system reset.

The Priority Inversion Problem

Priority order: $T1 > T2 > T3$

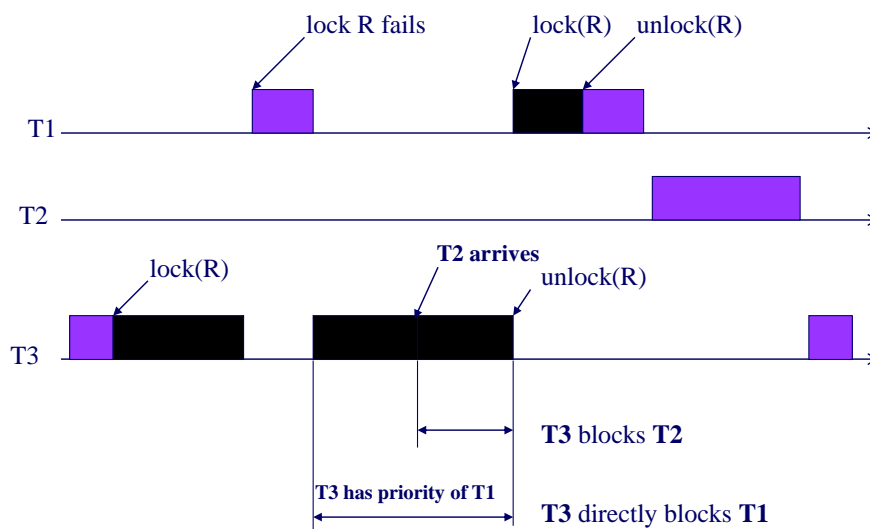


T2 is causing a higher priority task T1 wait !

Priority Inversion

1. T1 has highest priority, T2 next, and T3 lowest
2. T3 comes first, starts executing, and acquires some resource (say, a lock).
3. T1 comes next, interrupts T3 as T1 has higher priority
4. But T1 needs the resource locked by T3, so T1 gets blocked
5. T3 resumes execution (this scenario is still acceptable so far)
6. T2 arrives, and interrupts T3 as T2 has higher priority than T3, and T2 executes till completion
7. In effect, even though T1 has priority than T2, and arrived earlier than T2, T2 delayed execution of T1
8. This is "priority inversion" !! Not acceptable.
9. Solution T3 should inherit T1's priority at step 5

Priority Inheritance Protocol

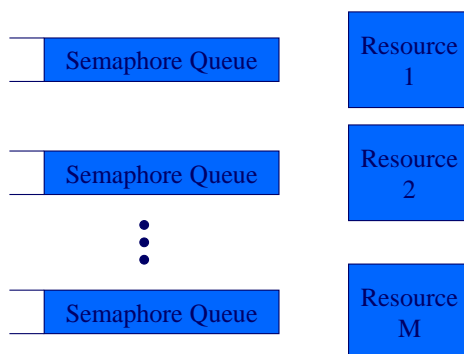


Priority Inheritance Protocol

- Question: What is the longest time a task can wait for lower-priority tasks?
- Answer: ?

Computing the Maximum Priority Inversion Time

- Consider the instant when a high-priority task that arrives.
 - What is the most it can wait for lower priority ones?



If I am a task, priority inversion occurs when

- (a) Lower priority task holds a resource I need (direct blocking)
- (b) Lower priority task inherits a higher priority than me because it holds a resource the higher-priority task needs (push-through blocking)

Schedulability Test

$$\forall i, 1 \leq i \leq n,$$

$$\frac{B_i}{P_i} + \sum_{k=1}^i \frac{C_k}{P_k} \leq i(2^{1/i} - 1)$$

Schedulability Test

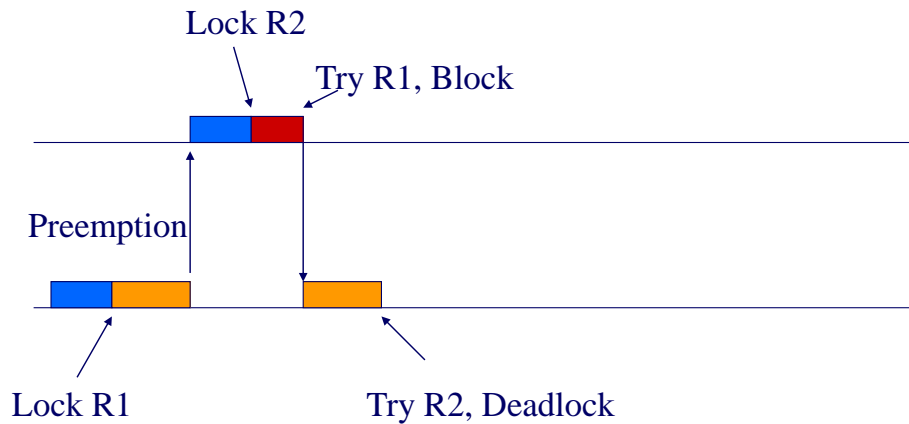
$$\forall i, 1 \leq i \leq n,$$

$$\frac{B_i}{T_i} + \sum_{k=1}^i \frac{C_k}{P_k} \leq i(2^{1/i} - 1)$$

Why do we have to test each task separately? Why not just one utilization-based test like it used to?

Problem: Deadlock

Deadlock occurs if two tasks locked two semaphores in opposite order



Spring '10

CIS 541

41

Priority Ceiling Protocol

- Definition: The priority ceiling of a semaphore is the highest priority of any task that can lock it
- A task that requests a lock R_k is denied if its priority is not higher than the highest priority ceiling of all currently locked semaphores (say it belongs to semaphore R_h)
 - The task is said to be blocked by the task holding lock R_h
- A task inherits the priority of the top higher-priority task it is blocking

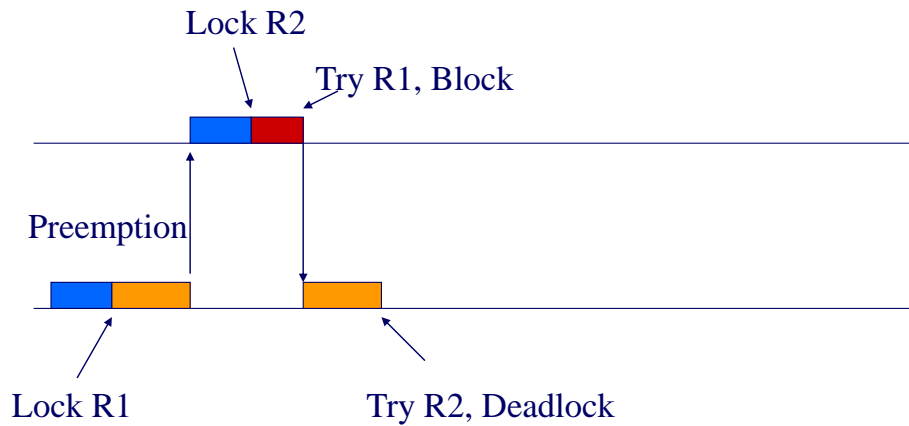
Spring '10

CIS 541

42

Problem: Deadlock?

Deadlock used to occur if two tasks locked two semaphores in opposite order. Can it still occur in priority ceiling?



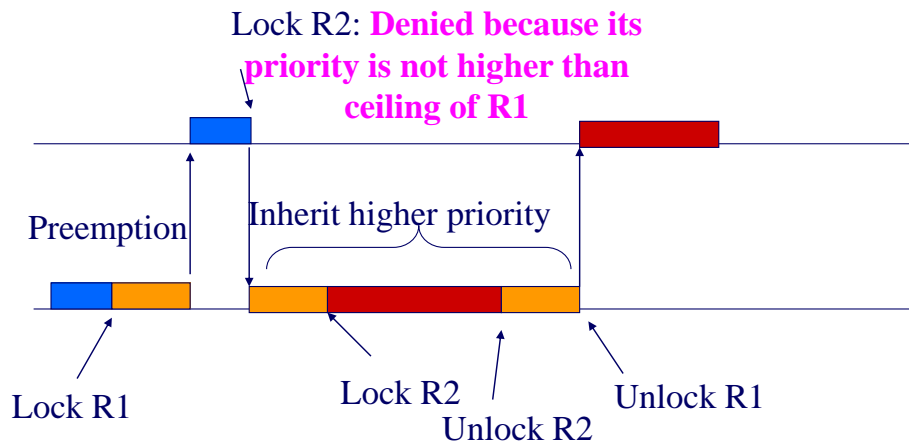
Spring '10

CIS 541

43

Problem: Deadlock?

Deadlock used to occur if two tasks locked two semaphores in opposite order. Can it still occur in priority ceiling?



Spring '10

CIS 541

44

Slack Resource Policy

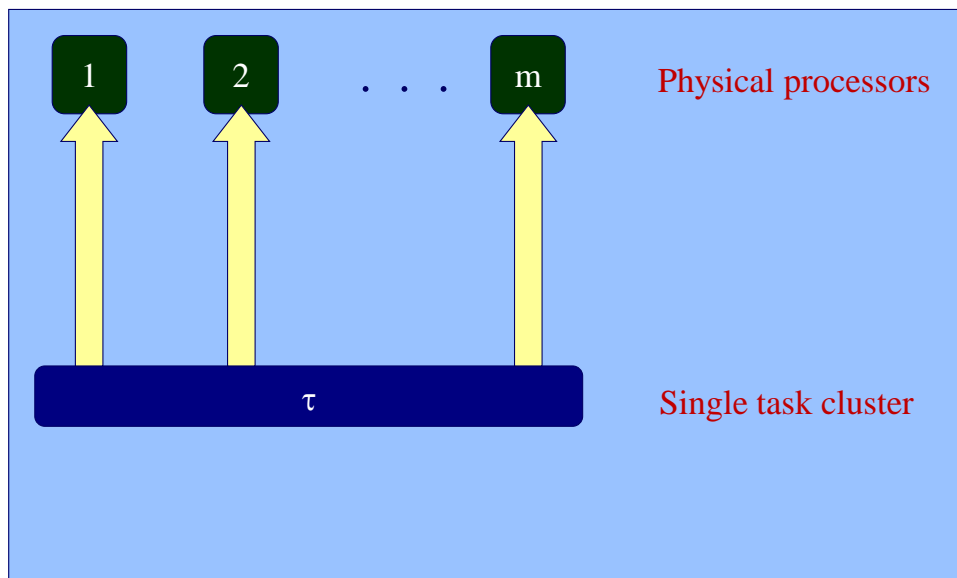
- **Priority:**
 - Any static or dynamic policy (e.g., EDF, RM, ...)
- **Preemption Level**
 - Any *fixed value* that satisfies: If A arrives after B and Priority (A) > Priority (B) then PreemptionLevel (A) > PreemptionLevel (B)
- **Resource Ceiling**
 - Highest preemption level of all tasks that might access the resource
- **System Ceiling**
 - Highest resource ceiling of all currently locked resources
- **A task can preempt another if:**
 - It has the highest priority
 - Its preemption level is higher than the system ceiling

MULTI-PROCESSOR SCHEDULING

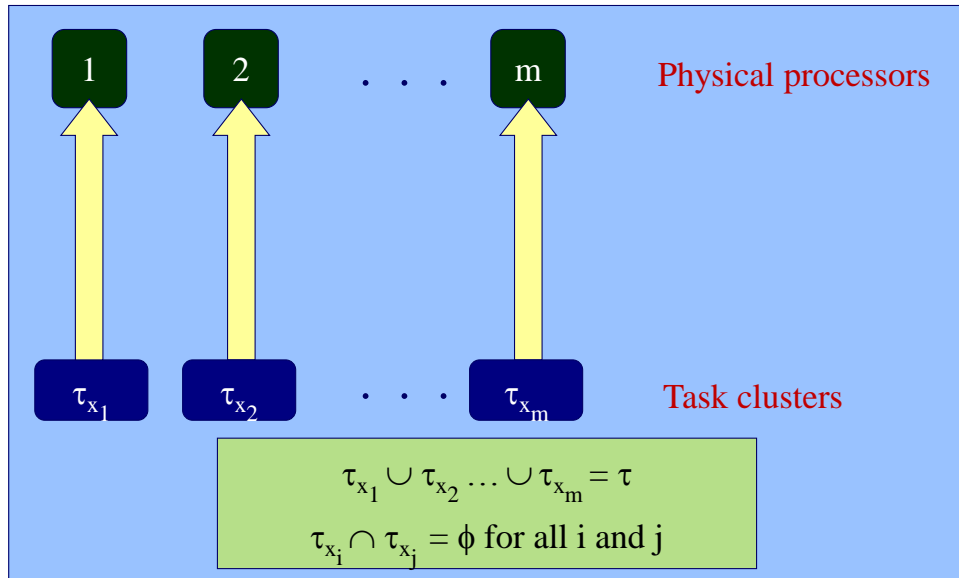
Multiprocessor Scheduling

- Why consider multiprocessors
 - Better tradeoff between computational power and costs (energy, fabrication)
 - Ability to exploit inherent concurrency in software
- Problem statement
 - Constrained deadline sporadic task system
 - $W = \{\tau_1, \dots, \tau_n\}$, where $\tau_i = (T_i, C_i, D_i)$ and $C_i \leq D_i \leq T_i$
 - C_i units must be supplied non-concurrently
 - Identical, unit-capacity multiprocessor platform
 - m processors
 - How can W be scheduled on these m processors?

Global Scheduling

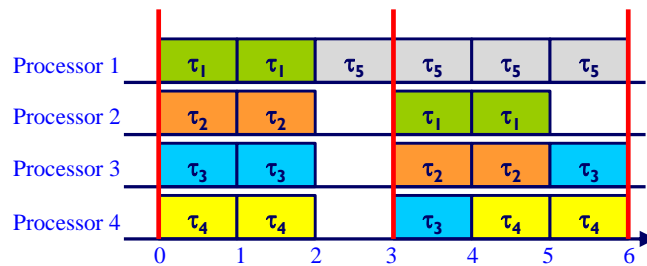


Partitioned Scheduling



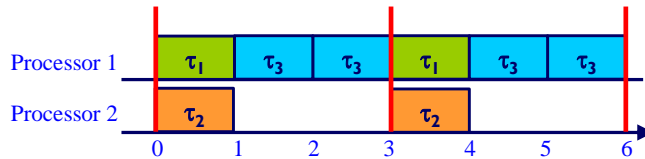
Counter-example for Partitioned

- Task set and number of processors
 - $\tau_1 = \tau_2 = \tau_3 = \tau_4 = (3, 2, 3)$ and $\tau_5 = (6, 4, 6)$, $m = 4$
- No partitioning technique will work
 - Some processor must be assigned two tasks which is not possible
- Schedulable under global LLF



Counter-example for Global

- Task set and number of processors
 - $\tau_1 = \tau_2 = (3, 1, 3)$ and $\tau_3 = (7, 6, 7)$, $m = 2$
- Global EDF cannot schedule this task set

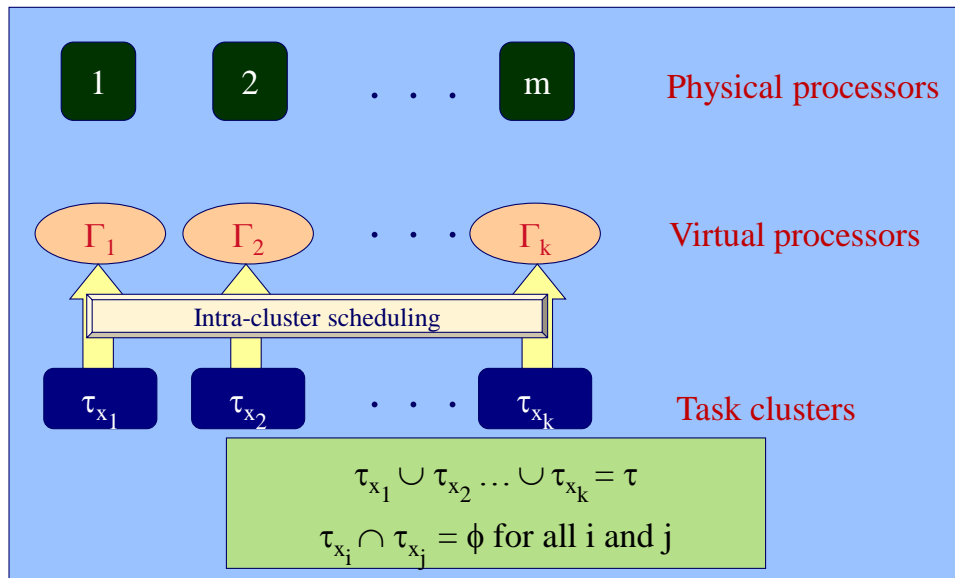


- Partitioned EDF can
 - τ_3 on processor 1, and the other two on processor 2

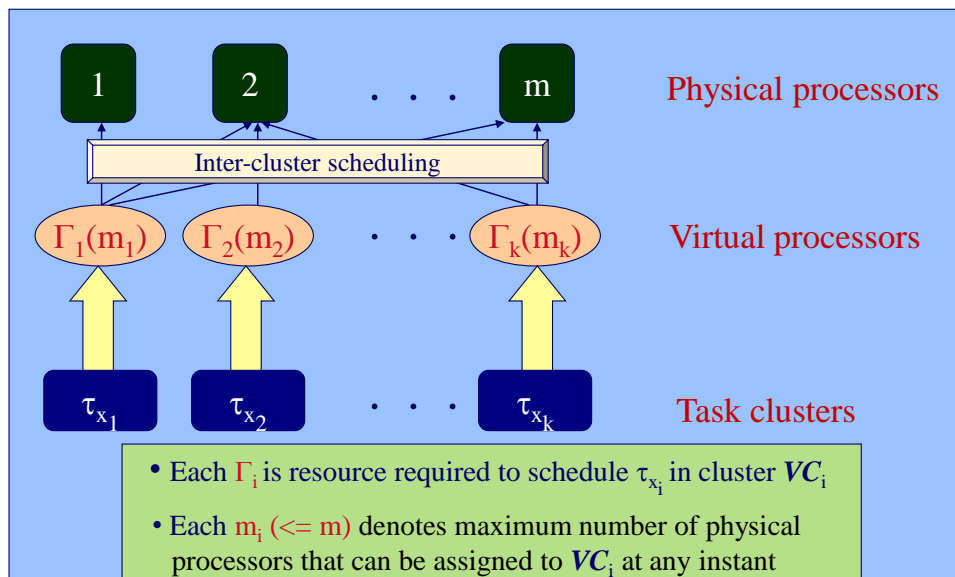
Partitioned vs. Global

- Two extreme cases of task-processor clustering
 - one-one (partitioned) vs. one-all (global)
- Both have advantages and disadvantages
 - **Partitioned**: Low preemptions, but low resource utilization
 - **Global**: High resource utilization with high preemptions
- Optimal scheduling on identical platforms
 - Only developed for implicit deadline task systems ($D_i = T_i$ for all i)
 - All known optimal schedulers are global strategies (Pfair [BCPV96])
 - **Problem open for constrained deadline periodic task systems**
 - Shown to be impossible for sporadic task systems
- **Can we support general task-processor clustering through the concept of platform virtualization (hierarchical scheduling)?**

Virtual Cluster-based Scheduling



Virtual Cluster-based Scheduling

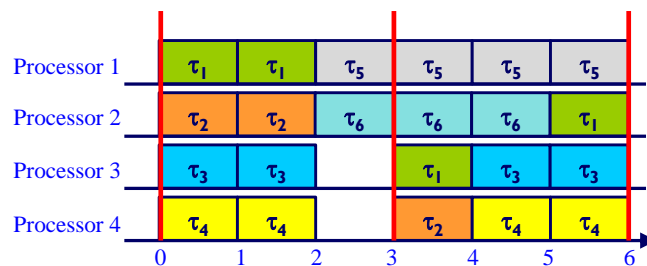


Counter-example for Partitioned and Global

- Task set and number of processors
 - $\tau_1=\tau_2=\tau_3=\tau_4=(3,2,3)$, $\tau_5=(6,4,6)$, and $\tau_6=(6,3,6)$, $m=4$
- No Partitioning technique can work
 - Some processor needs to be assigned two tasks which is not possible (maximum utilization that can be assigned to any processor is 1)

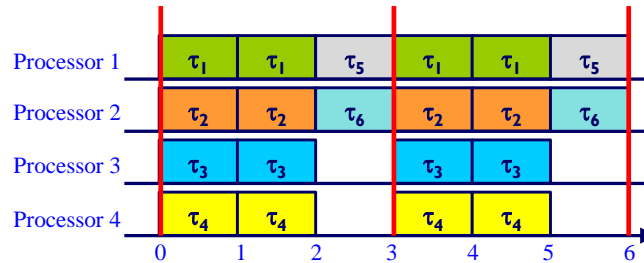
Counter-example for Partitioned and Global

- Task set and number of processors
 - $\tau_1=\tau_2=\tau_3=\tau_4=(3,2,3)$, $\tau_5=(6,4,6)$, and $\tau_6=(6,3,6)$, $m=4$
- Schedule under global EDF/EDZL (earliest deadline until zero laxity)/LLF
 - Task τ_2 misses its deadline



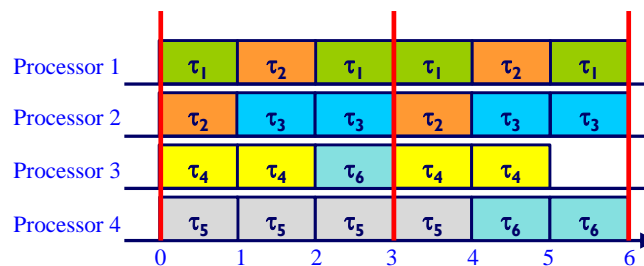
Counter-example for Partitioned and Global

- Task set and number of processors
 - $\tau_1=\tau_2=\tau_3=\tau_4=(3,2,3)$, $\tau_5=(6,4,6)$, and $\tau_6=(6,3,6)$, $m=4$
- Schedule under global fp-EDF/US-EDF (highest priority to high utilization tasks τ_1, \dots, τ_5)
 - Task τ_6 misses its deadline



Virtual Clustering

- Task set and number of processors
 - $\tau_1=\tau_2=\tau_3=\tau_4=(3,2,3)$, $\tau_5=(6,4,6)$, and $\tau_6=(6,3,6)$, $m=4$
- Schedule under clustered scheduling
 - τ_1, τ_2, τ_3 scheduled on processors 1 and 2
 - τ_4, τ_5, τ_6 scheduled on processors 3 and 4



Virtual Clustering

- Two-level hierarchical scheduler
 - Intra-cluster schedulers for tasks within clusters
 - Inter-cluster schedulers for clusters on the platform (clusters can share some physical processors)
- Concurrency bound for each cluster
 - Abstract concurrency constraints of tasks within cluster
 - Helps regulate resource access (e.g., Caches)
- Have virtual clusters been used before?
 - Supertasks^[MoRa99], Megatasks^[ACD06]
 - Results restricted to Pfair schedulers (not generalizable)

MODE CHANGE PROTOCOLS