# Task Synchronization on Multiprocessors

## Raj Rajkumar
Lecture #18

---

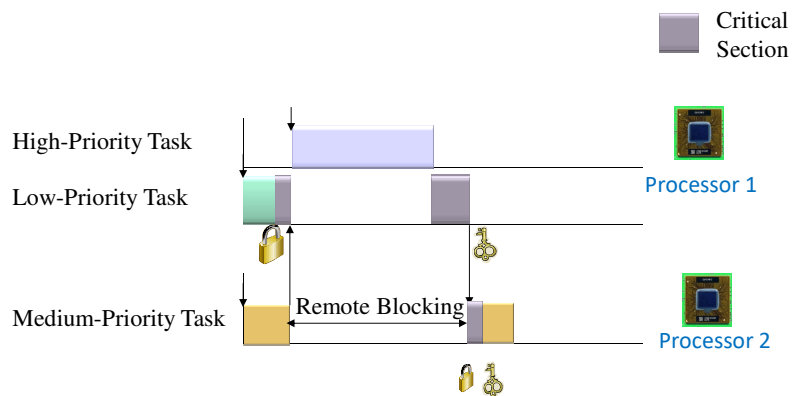# Multiprocessor Task Synchronization

- Multiprocessor synchronization challenges
  - Remote blocking
  - Self-suspending behavior
  - Multiple priority inversions due to suspensions

- Uniprocessor solutions need to be revisited
  - Priority inheritance protocols
  - Priority ceiling protocol

## Revisiting Blocking from Resource Sharing

- Consider the blocking term $B_i$ of a task $\tau_i$ as the **additional time penalty** incurred by $\tau_i$ due to resource sharing **relative to the case when there is absolutely _no_ resource sharing** in the taskset.
- On a uniprocessor, when one of the priority inheritance protocols is used, only the critical sections of lower priority tasks can contribute to $B_i$
  - Under the **basic priority inheritance protocol**, up to *(m, n)* critical sections of lower priority critical sections, where *n* is the number of tasks with lower priority than $\tau_i$ and *m* is the number of mutexes accessed by lower-priority tasks with a priority ceiling higher than or equal to the priority of $\tau_i$
  - Under the **priority ceiling protocol**, at most one critical section of a lower-priority task with a priority ceiling higher than or equal to the priority of $\tau_i$
  - The **highest locker protocol** has the same property as the priority ceiling protocol
  - Under the **non-preemption protocol**, at most any critical section of a lower-priority task

---

## Remote Blocking

- Remote Blocking
  - Tasks wait on resources held on other cores
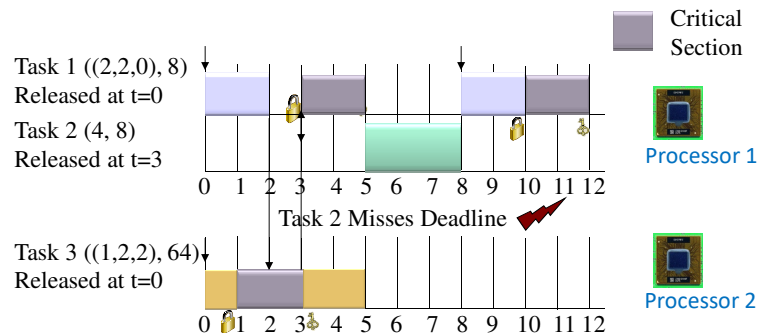  - Leads to wastage of CPU utilization

2

# Multiprocessor Resource Sharing

- If a task shares a resource with a task on any other (remote) processor, it may pay an additional penalty.
- So, sharing resources with tasks of **any** priority on **any** other processor can contribute to an increase in $B_i$.
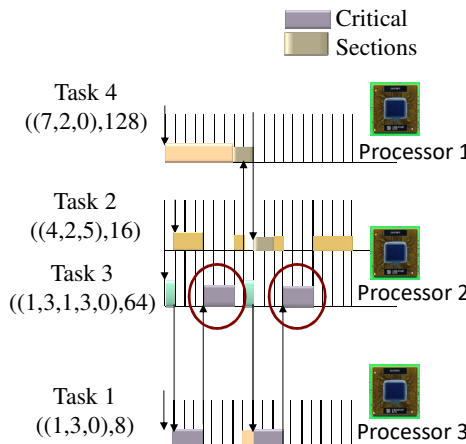
---

# Self-Suspending Tasks

- Tasks suspend on remote cores
  - In the uni-processor context
    - Tasks suspend on other tasks in the same processor
  - Leads to various scheduling penalties



Task 1 ((2,2,0), 8)
Released at t=0

Task 2 (4, 8)
Released at t=3

0  1  2  3  4  5  6  7  8  9  10  11  12

Task 2 Misses Deadline

Task 3 ((1,2,2), 64)
Released at t=0

0  1  2  3  4  5  6  7  8  9  10  11  12

Critical Section

Processor 1

Processor 2

# Multiple Priority Inversions

- When a task suspends
  - Lower-priority tasks could be released
  - Lower-priority tasks could acquire resources
    - Due to priority ceiling, they can cause preemptions
- The worst case:
  - For each task segment,
  - One preemption from each lower-priority task

Critical Sections

Task 4
((7,2,0),128)

Processor 1

Task 2
((4,2,5),16)

Processor 2

Task 3
((1,3,1,3,0),64)

Task 1
((1,3,0),8)

Processor 3

**Carnegie Mellon**

*18-648: Embedded Real-Time Systems*
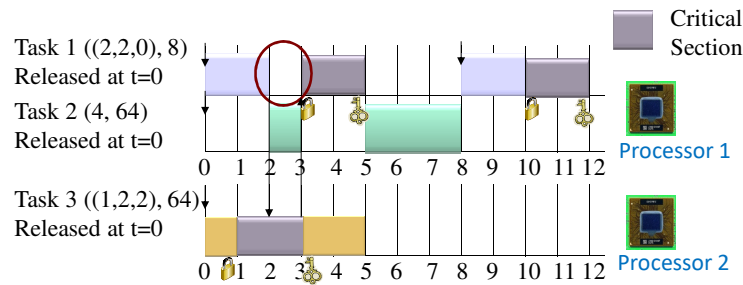
Electrical & Computer
ENGINEERING

---

# Priority Ceiling Protocol Extension

- Multiprocessor Priority Ceiling Protocol (MPCP)
  - **Global critical sections** guarded by **global mutexes**
  - Global mutexes associated with global priorities
    - **(Global) Priority ceiling of global mutex $G_M$** is calculated as:

    > Max (All task normal execution priorities) +
    >     1 + Max (Normal priorities of tasks accessing $G_M$)

- Critical section executions can be preempted
  - By other critical section executions with an even greater priority ceiling
- Each mutex $M_G$ has a priority queue
  - Pending tasks acquire mutex in priority order

**Carnegie Mellon**

*18-648: Embedded Real-Time Systems*

Electrical & Computer
ENGINEERING

4

# Suspension-based Protocol

- When global mutex request is pending
  - Task suspended and other active tasks execute
- Issues:
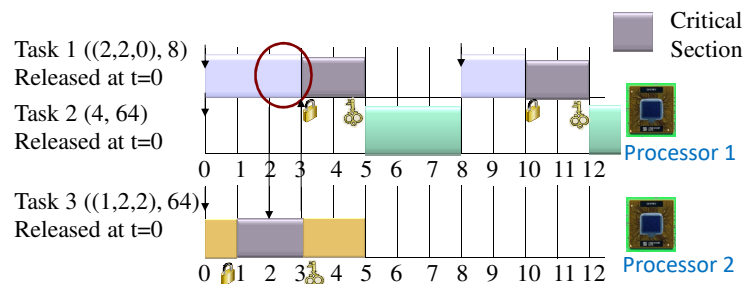  - Self-suspension, multiple priority inversions

Task 1 ((2,2,0), 8)
Released at t=0

Task 2 (4, 64)
Released at t=0

0 1 2 3 4 5 6 7 8 9 10 11 12

Processor 1

Task 3 ((1,2,2), 64)
Released at t=0

0 1 2 3 4 5 6 7 8 9 10 11 12

Processor 2

Critical Section

---

# Spinning-based Protocol

- When a global mutex request is pending
  - Task continues to "spin" till it acquires resource
- Issues:
  - Spinning-based utilization loss

Task 1 ((2,2,0), 8)
Released at t=0

Task 2 (4, 64)
Released at t=0

0 1 2 3 4 5 6 7 8 9 10 11 12

Processor 1

Task 3 ((1,2,2), 64)
Released at t=0

0 1 2 3 4 5 6 7 8 9 10 11 12

Processor 2

Critical Section
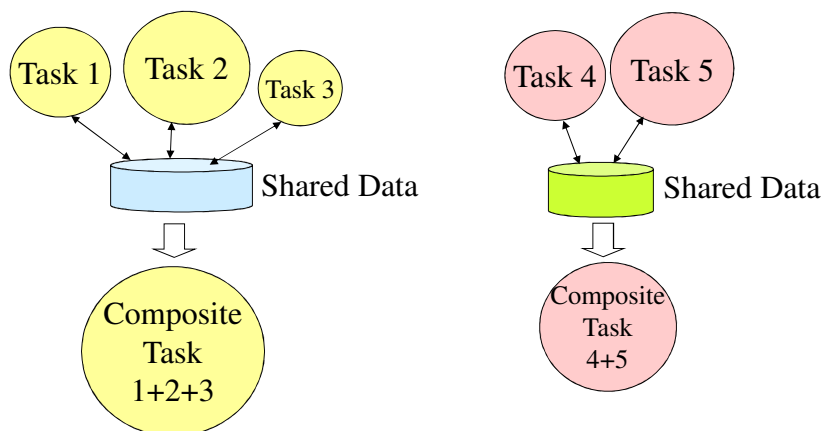
## Evaluating The Effectiveness of Bin-Packing

- Suppose the number of bins required by scheme $S$ to pack a set of objects is $|S|$.
- One may want to compare the performance of a bin-packing heuristic $H$ with that of the *optimal* scheme $O$
- Analyses:
  - **Worst-case analysis**: theoretical study to determine the absolute worst-case number of bins required by $H$ relative to that of $O$
    - i.e. find $max(|H| \div |O|)$
  - **Average-case analysis**: compute $(|H| \div |O|)$ for randomly picked sets of objects
    - Perform over a large sample so that results are statistically significant
    - How do we know $|O|$?
      - Finding the optimal packing is NP-hard:
      - Will take an exponential amount of time (unless P = NP)
    - Create by construction!
      - Assume fully packed bins: split into known # of objects in each bin into random sizes (for a given # of bins = $|O|$)
      - Unpack these objects and offer to heuristic for packing
      - Original packing we started with is optimal!

---

# Synchronization-Aware Bin-Packing

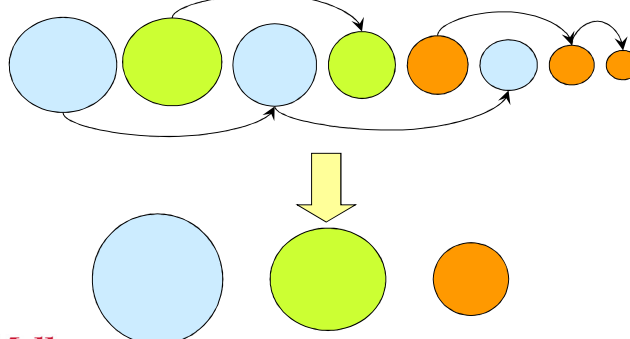# Making Bin-Packing Synchronization-Aware

- Global critical sections lead to
  - Significantly more blocking (i.e. the $B_i$ term becomes much bigger)
    - "Blocking" can be caused both by-higher priority <u>and</u> lower-priority tasks on other processors.
  - Deferred execution behaviors can result in scheduling penalties
- Conversely, local critical sections
  - Utilize the (local) priority ceiling protocols (like the highest locker priority) to bound $B_i$ to a single (local) critical section of a lower priority task
  - Experience no deferred execution penalties
- Our bin-packing heuristics use object size as the primary criterion while packing objects into bins
- Can we try to exploit the bin-packing heuristics and try to make "global critical sections" become local ones?
  - Combine objects (tasks) into **composite objects**.

**Carnegie Mellon**     *18-648: Embedded Real-Time Systems*     Electrical **&** Computer ENGINEERING

---

# Synchronization-Aware Packing



**Carnegie Mellon**     *18-648: Embedded Real-Time Systems*     Electrical **&** Computer ENGINEERING

7

# Bin-Packing → Packing of Composite Objects

- If tasks A and B share data, they become a composite task
  - Applies recursively: task A and/or B can also be a composite task
  - The relationship is transitive:
    - i.e. if tasks A and B share one piece of data, and tasks B and C share another piece of data, tasks A, B and C become a single composite task
- Transform the given list of tasks with their shared resource accesses into a shorter list of bigger composite tasks

---

# Bin-Packing of Composite Objects

- Use bin-packing heuristics to allocate composite tasks to bins
  1. If a composite task does not fit into a bin, set it aside and allocate other unallocated tasks
  2. After current iteration of allocation, if any unallocated tasks remain, then try to allocate them
     - This reduces the number of global critical sections
  3. If a composite task does not fit any bin, split into two objects with *minimum synchronization cost* and allocate. If a task cannot be split (i.e. it is an *elementary* task and not a *composite* task), you must create a new bin.
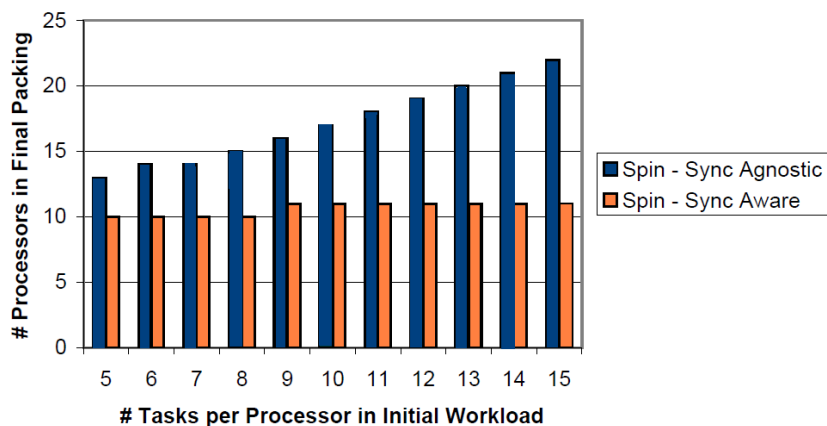  4. Go to Step 1.

  Other variations are possible.
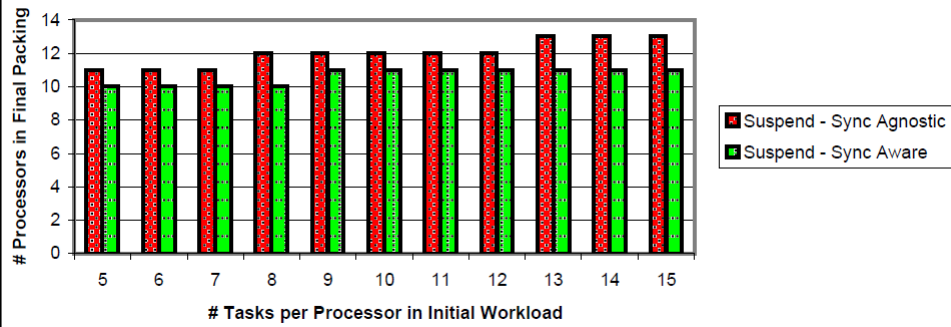
## Comparison of Schemes

- Synchronization itself between the objects <u>after</u> allocation can be either spin-based or suspension-based.
- Evaluation is *very* dependent on assumptions:
  - Distribution of task utilization (object sizes)
  - Amount of sharing (# of tasks that share the same resources)
  - Number of resource accesses per task
  - Length of each critical section
  - Overhead of preemptions
  - Overhead of caching penalty during preemptions

## Synchronization-Aware Bin Packing - Spinning



8 fully packed cores, 2 critical sections per task, 2 lockers per mutex
500μs critical sections, Periods: [10ms, 100ms]

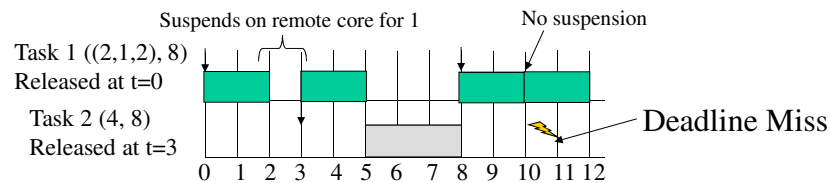## Synchronization-Aware Bin Packing - Suspension



8 fully packed cores, 2 critical sections per task, 2 lockers per mutex
500µs critical sections, Periods: [10ms, 100ms]
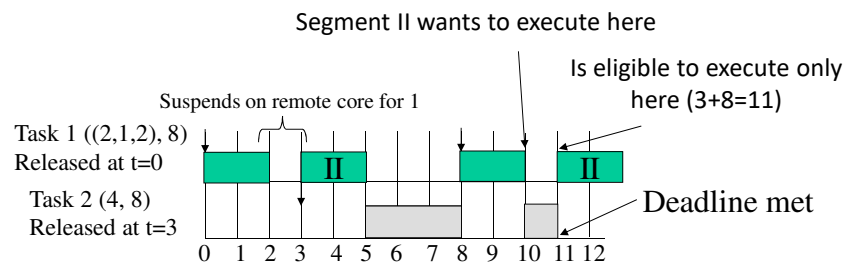
---

# Dealing with Deferred Execution

## Self-Suspending Tasks

- Tasks can suspend due to variety of reasons:
  - Sensor Input
  - Synchronization with other processing core
  - Waiting on co-processor
- Consider the example:
  - Task 1: Period 8: Executes for 2 units, Suspends for at most 1 unit, Continues execution for 2 more units
  - Task 2: Period 8: Execution for 4 units

## Solution: Period Enforcement

- An execution segment becomes eligible to execute only at periodic boundaries.
- If a task suspends itself during execution (for any reason including suspension on a shared resource that is locked by a task on a remote processor), execution before and after suspension is each considered an execution segment.

## Conclusions

- Both bin-packing and partitioned scheduling can be improved by "splitting" tasks
  - At most one task needs to be split per processor
- Period transformation can be used to reconcile scheduling priority and semantic importance of tasks
- With harmonic tasks and splitting, we can get 100% schedulable utilization per processor
- Multiprocessor synchronization schemes
  - Period enforcement
  - Static and dynamic slack enforcement