

## 7.c Global resource sharing

### Multiprocessor PCP /1

- P-FPS with resources bound to processors [Sha, Rajkumar, Lehoczky, 1988]
  - The processor that hosts a resource is the *synchronization processor* (SP) for that resource
    - It statically knows all the use requirements of all of its resources
  - The critical sections of a resource execute on its SP
    - Jobs that use *remote* resources employ “*distributed transactions*”
  - The processor to which a task is assigned is the *local processor* (LP) for all of the jobs of that task

### Contention and blocking




- The single-runner premise on which previous solutions were based falls apart
  - Suspending on wait no longer facilitates earlier release of shared resources ← parallelism gets in the way
  - Priority boosting the lock holder does not help ← per-CPU priorities do not have global meaning (on partitioned scheduling)
  - With local *and* global resources, suspensive wait becomes dangerous ← local priority inversions (PI) may occur
  - Spinning protects against PI, but wastes CPU cycles

### Multiprocessor PCP /2

- A task may use local *and* global resources
  - Local resources reside on the LP of that task
  - Resources are global when their SP differs from the client tasks' LP
- Resource access control protocols need *actual locks* to protect against parallel contention
  - This causes *lock-free algorithms* to become attractive
- SPs use M-PCP to control access to their global resources


## Multiprocessor PCP /3

- The task that holds a global lock should *not* be preempted locally
  - All global critical sections must execute at higher ceiling priorities than all local tasks on their SP
  - **This breaks independence!** 
- A task  $\tau_h$  that is denied access to a global shared resource  $\rho_g$  suspends on its LP and waits in a priority-based queue for that resource
  - Any task  $\tau_l$  with lower-priority than  $\tau_h$  on the same LP may thus acquire global resources on  $\rho_g$ 's SP, with higher ceiling than  $\rho_g$

## Blocking under M-PCP

- With M-PCP, task  $\tau_i$  is *blocked* by lower-priority tasks in 5 ways!
  - *Local blocking (once per release)*: when finding a local resource held by a local lower-priority task that got running as a consequence of  $\tau_i$ 's suspension on access to a global resource
  - *Remote blocking (once per request)*: when finding a global resource held by a lower-priority task running on the global resource's SP
  - *Local preemption*: when global critical sections are executed on  $\tau_i$ 's LP by remote tasks of any priority (multiple times) and by local tasks of lower priority (once per release)
  - *Remote preemption (once per request)*: when higher-ceiling global critical sections execute on the SP where  $\tau_i$ 's global resource resides
  - *Deferred interference* as local higher-priority tasks suspend on access to global resources because of blocking effects

## Multiprocessor PCP /4

- If the global resource  $\rho_{g'}$  being acquired by  $\tau_l$  resides on the same SP as  $\rho_g$ , then  $\tau_h$  suffers an anomalous form of PI
  - The execution in  $\rho_{g'}$  delays the release of  $\rho_g$
- As contention for a global resource involves suspension, **M-PCP suffers the risk of deadlock** 
  - With global resources hosted on  $> 1$  SPs, nesting global resources may lead to deadlock and *must be disallowed*
- This is why other protocols prefer  $\tau_h$  to spin

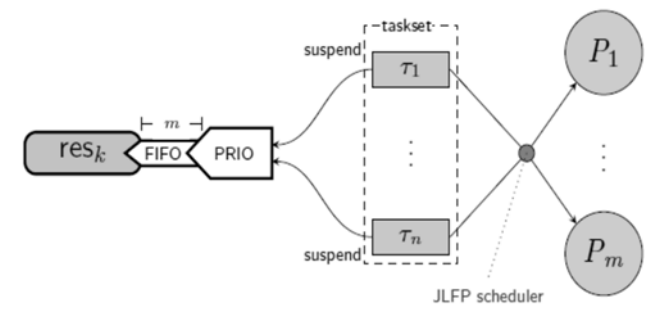
## Multiprocessor SRP

- P-EDF with resources bound to processors [Gai, Lipari, Di Natale, 2001]
  - Normal SRP is used for controlling access to local resources
- Tasks that lock a global resource execute the critical section at the highest local priority
  - As the lock-holder cannot be pre-empted, the wait time is shorter
  - But this provision breaks independence
- Tasks that request a global resource  $\rho_G$  already locked, are held in a FIFO queue on  $\rho_G$ 's SP and *spin* on their LP
  - This policy upper-bounds the requesting task's *spin time* to  $m - 1$  executions of the longest critical section of  $\rho_G$
  - This duration adds to the task's WCET

## In general ...

- With lock-based resource access control protocols, locks can use either *suspension* or *spinning*
- With suspension, the calling task that cannot acquire the lock is placed in a priority-ordered queue
  - To bound blocking, PI avoidance algorithms should be used
- With spinning, the task busy-waits
  - To bound blocking, the spinning task becomes non-preemptible and its request is placed in FIFO queue
- The lock owner may also run non-preemptively
  - But this breaks independence

## $O(m)$ locking protocols : G-EDF /2



## $O(m)$ locking protocols : G-EDF /1

- All resources are global
  - To request a resource, a task must first acquire a general priority-queue, PQ, lock (one of  $m$ )
  - If the resource is busy, the requestor waits, *suspending*, on a resource-specific FIFO queue, FQ (of  $m$  positions)
  - The lock-holder inherits the highest priority of tasks waiting in the chain of queues (FQ and PQ)
- Per-request blocking is  $2m - 1$  executions of the longest critical section for the resource
  - When FQ is full with  $m$  lp-jobs and  $m$  hp-tasks run (including the job of interest) that all want to access the same resource
- The other tasks suffer inheritance blocking

## $O(m)$ locking protocols : P-EDF /1

- Shared resources may be local or global
  - One *priority queue* (PQ) per processor: the task at the head of it acquires a *token* to use to contend for global resources
  - Requests for G-resources wait in a per-resource FQ
    - The waiting tasks suspend
  - Lock-holders' priority is inheritance-boosted from their PQ
- Blocking for all tasks has three components
  - *Local*, when the lock-holder is a local lp-task (per release)
  - *Remote direct*, when the requestor is last in the FQ (per request)
  - *Remote transitive*, when a local lp-task has acquired the PQ token and is last of the FQ (per release)

### $O(m)$ locking protocols : P-EDF /2

binary semaphore and prio boosting

2017/18 UniPD - T. Vardanega Real-Time Systems 524 of 595

### $O(m)$ independence preservation /2

- Clusters of size  $1 \leq c \leq m$ 
  - Global scheduling per cluster, partitioned cluster assignment
- *Suspension-based*
  - One FIFO+PRIO queue per cluster, for  $O(m)$  blocking
  - One per-resource global FIFO queue
    - Head of cluster FQ copied in G-FQ and removed only after service
- Independence preserved by *inter-cluster migration*
  - Head of G-FQ (if pre-empted) can migrate to any CPU along the queue (hence across clusters), with priority boosted by inheritance from a waiting task
- Blocking is *per request*:  $\beta_{i,k} = (m - 1)\omega_k$

2017/18 UniPD - T. Vardanega Real-Time Systems 526 of 595

### $O(m)$ independence preservation /1

JLFP scheduler

$$\beta_{i,k} = \left[ \left( \frac{m}{c} - 1 \right) + \frac{m}{c} \times (c - 1) \right] \omega_k = (m - 1)\omega_k$$

2017/18 UniPD - T. Vardanega Real-Time Systems 525 of 595

### $O(m)$ independence preservation /3

- $t = 3$ : task  $\tau_2$  suspends and task  $\tau_1$  resumes execution
- $t = 4$ : task  $\tau_3$  migrates to cluster<sub>1</sub> and preempts task  $\tau_1$

2017/18 UniPD - T. Vardanega Real-Time Systems 527 of 595

## [Brandenburg, 2013]

## ■ Theorem

- Under non-global scheduling (with cluster size  $c < m$ ), *no* resource access control protocol can simultaneously
  - Prevent unbounded PI blocking
  - Preserve independence (you don't suffer if you don't contend)
  - Avoid migration
- *Seeking independence preservation and bounded PI-blocking requires inter-cluster job migration (!)*

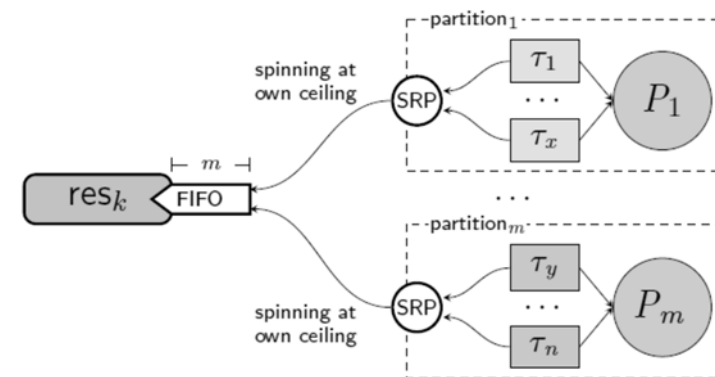
## MrsP [Burns, Wellings, 2013] /2

- Spinning non-preemptively may decrease feasibility
  - Urgent tasks would suffer longer blocking
- Spinning at the *local* ceiling priority is better
  - With all processors using PCP/SRP, *at most one task per processor* may contend globally, which assures  $O(m)$  blocking
  - Access requests are served in FIFO order
- To bound blocking, spinning tasks “donate” their cycles to the pre-empted lock-holder
  - The lock-holder migrates to the processor of a spinning task and runs in its stead until it either completes or migrates again

## MrsP [Burns, Wellings, 2013] /1

- Rendering RTA for partitioned multiprocessors *identical* to the single-processor case
  - The cost of accessing global resources should be *increased* to reflect the need to serialize parallel contention
- Preserving the property that, once a task starts executing, its resources *are* available
  - It needs global resource control protocols
  - Cannot use suspension-based solutions!

## MrsP [Burns, Wellings, 2013] /3



## MrsP [Burns, Wellings, 2013] /4

- For partitioned scheduling ( $c = 1$ )
- *Spinning-based*: local wait spins at local ceiling
  - Combined with PCP/SRP, this assures blocking at most once before execution
- Allows using uniprocessor-style RTA
- Wait is *per resource*, increased by parallel contention
  - $\beta_i = \max_k(\omega_k^{MrsP}) = \max_k((m-1)\omega_k) = (m-1) \times \max_k(\omega_k)$
- Earlier release obtained by migrating lock holder (if preempted) to the CPU where the first contender in the global FIFO is currently spinning

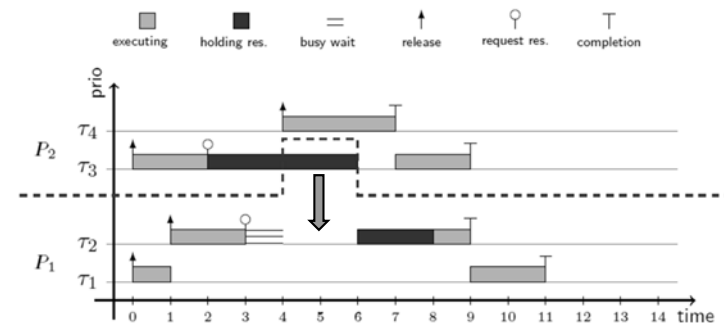
## MrsP [Burns, Wellings, 2013] /6

- Resource nesting can be supported with either *group locking* or *static ordering* of resources
  - With static ordering, resource access is allowed only with order number greater than any currently held resources
  - The implementation should provide an «out of order» exception to prevent run-time errors
- The ordering solution is better than banning nesting and has less penalty than group locking
- Recent work has extended MrsP to proper nesting

## MrsP [Burns, Wellings, 2013] /5

- $R_i = C'_i + B_i + I_i$
- $B_i = \max\{\rho_l, b\}$ 
  - $\rho_l$  is the longest critical section of a resource used by a lower-priority task with ceiling no less than  $\tau_i$ 's priority
  - $b$  is the longest duration of RTOS inhibited preemption
- $I_i = \sum_{j \in \text{hpl}(i)} \left\lceil \frac{R_i}{T_j} \right\rceil C'_j$
- $C'_i = C_i + \sum_j n_j e_j$ 
  - $C_i$  is task  $\tau_i$ 's WCET outside of critical sections
  - $n_j$  is the number of times task  $\tau_i$  uses shared resource  $j$
  - $e_j \leq (m-1)\rho_j$ , with  $\rho_j$  the longest critical section of resource  $j$

## MrsP [Burns, Wellings, 2013] /7



- $t = 3$ : task  $\tau_2$  start spinning at ceiling priority
- $t = 4$ : task  $\tau_3$  migrates to  $P_1$  and executes in place of  $\tau_2$

## Summary

- Issues and state of the art
- Dhall's effect: examples
- Scheduling anomalies: examples
- P-fair scheduling
- Sufficient tests for simple workload model
- Recent extensions: DP-Fair and RUN
- Incorporating global resource sharing