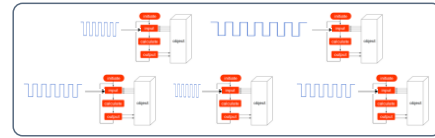


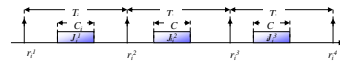
Fixed-Priority Multiprocessor Scheduling [RTAS 2010]

Joint work with
 Nan Guan, Martin Stigge and Yu Ge
 Northeastern University, China
 Uppsala University, Sweden

Real-time Systems



□ N periodic tasks (of different rates/periods)



Utilization/workload: C_i/T_i

□ How to schedule the jobs to avoid deadline miss?

On Single-processors

- Liu and Layland's Utilization Bound [1973] (the 19th most cited paper in computer science)

$$\sum_{\tau_i \in \tau} U_i \leq N(2^{1/N} - 1)$$

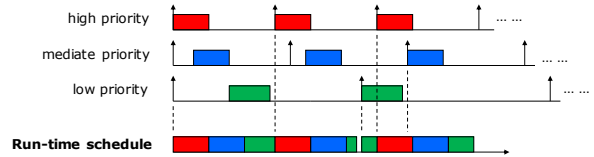
⇒ the task set is schedulable

number of tasks

- $N \rightarrow \infty, N(2^{1/N} - 1) = 69.3\%$
- Scheduled by RMS (Rate Monotonic Scheduling)

Rate Monotonic Scheduling

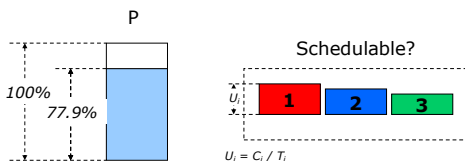
- Priority assignment: shorter period → higher prio.
- Run-time schedule: the highest priority first



□ How to check whether all deadlines are met?

Liu and Layland's Utilization Bound

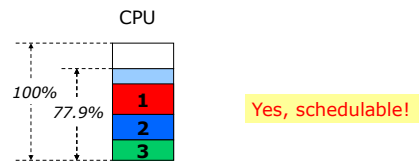
- Schedulability Analysis



Liu and Layland's bound:
 $3 \times (2^{1/3} - 1) = 77.9\%$

Liu and Layland's Utilization Bound

- Schedulability Analysis



Liu and Layland's bound:
 $3 \times (2^{1/3} - 1) = 77.9\%$

Multiprocessor (multicore) Scheduling

- ❑ Significantly more difficult:
 - Timing anomalies
 - Hard to identify the worst-case scenario
 - Bin-packing/NP-hard problems
 - Multiple resources e.g. caches, bandwidth
 -

Open Problem (since 1973)

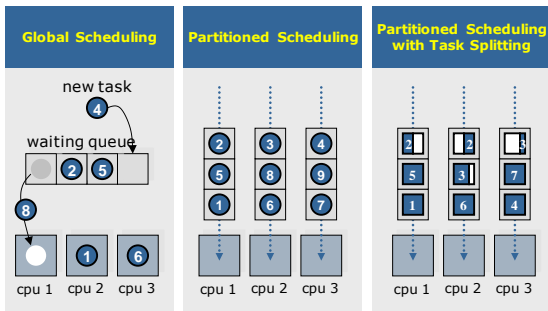
- ❑ Find a multiprocessor scheduling algorithm that can achieve Liu and Layland's utilization bound

$$\frac{\sum C_i/T_i}{M} \leq N(2^{1/N} - 1)$$

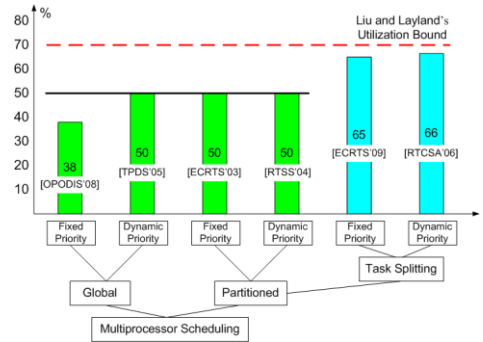
⇒ the task set is schedulable

number of processors

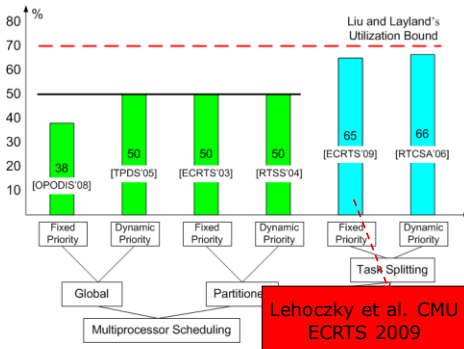
Multiprocessor Scheduling



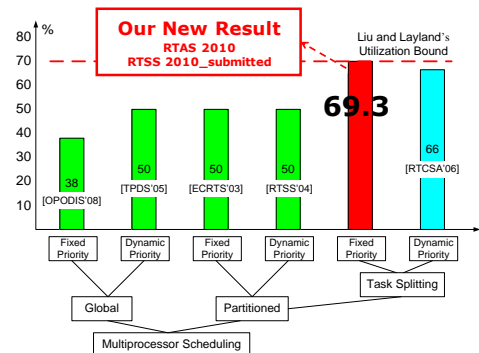
Best Known Results (before 2010)



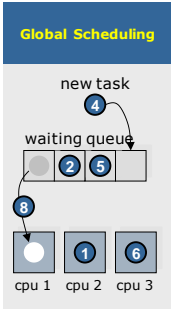
Best Known Results (before 2010)



Best Known Results

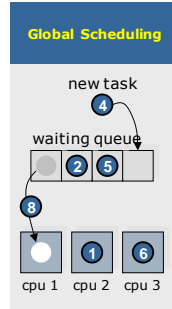


Multiprocessor Scheduling



Would fixed-priority scheduling e.g. "RMS" work?

Multiprocessor Scheduling

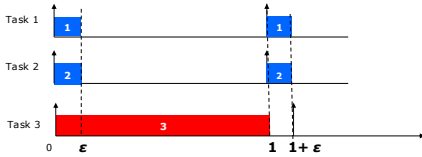


Would fixed-priority scheduling e.g. "RMS" work?

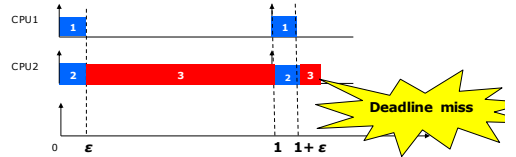
Unfortunately "RMS" suffers from the **Dhall's anomaly**

Utilization may be "0%"

Dhall's anomaly



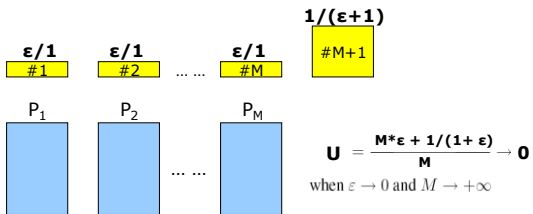
Dhall's anomaly



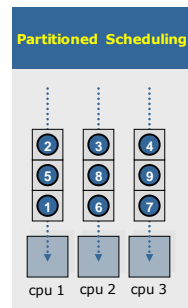
Schedule the 3 tasks on 2 CPUs using "RMS"

Dhall's anomaly

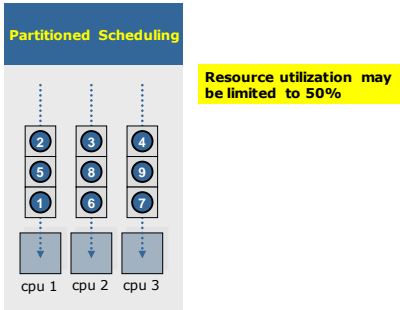
(M+1 tasks and M processors)



Multiprocessor Scheduling



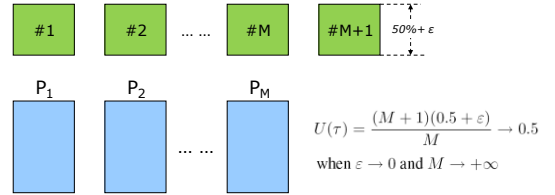
Multiprocessor Scheduling



Partitioned Scheduling

- The Partitioning Problem is similar to Bin-packing Problem (NP-hard)

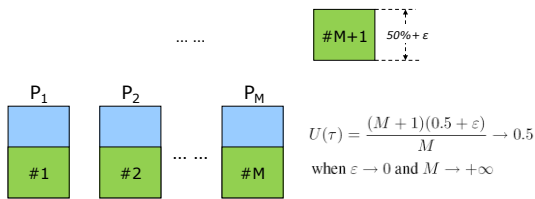
- Limited Resource Usage, 50% $\sum C_i/T_i \leq 1$
necessary condition to guarantee schedulability



Partitioned Scheduling

- The Partitioning Problem is similar to Bin-packing Problem (NP-hard)

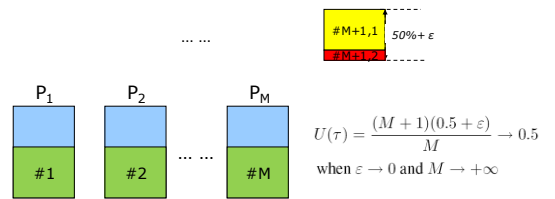
- Limited Resource Usage $\sum C_i/T_i \leq 1$
necessary condition to guarantee schedulability



Partitioned Scheduling

- The Partitioning Problem is similar to Bin-packing Problem (NP-hard)

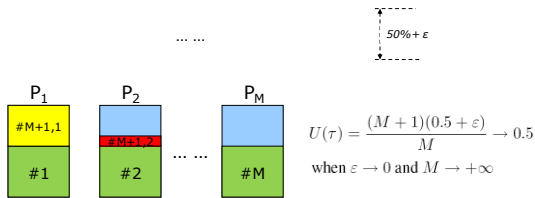
- Limited Resource Usage $\sum C_i/T_i \leq 1$
necessary condition to guarantee schedulability



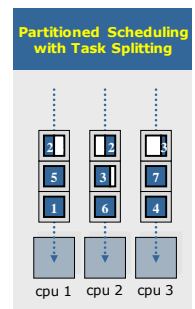
Partitioned Scheduling

- The Partitioning Problem is similar to Bin-packing Problem (NP-hard)

- Limited Resource Usage $\sum C_i/T_i \leq 1$
necessary condition to guarantee schedulability

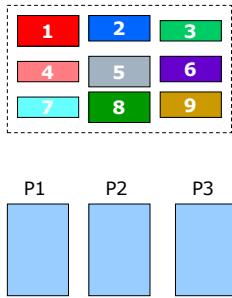


Multiprocessor Scheduling



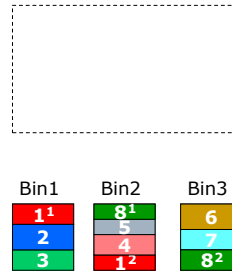
Partitioned Scheduling

- Partitioning



Bin-Packing with Item Splitting

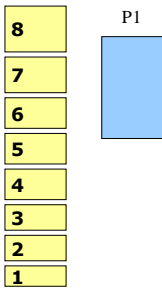
- Resource can be "fully" (better) utilized



Previous Algorithms

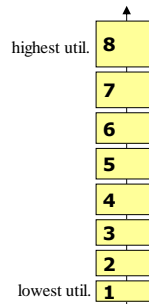
[Kato et al. IPDPS'08] [Kato et al. RTAS'09] [Lakshmanan et al. ECRTS'09]

- Sort the tasks in some order e.g. utilization or priority order
- Select a processor, and assign as many tasks as possible



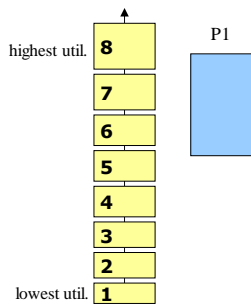
Lakshmanan's Algorithm [ECRTS'09]

- Sort all tasks in decreasing order of utilization



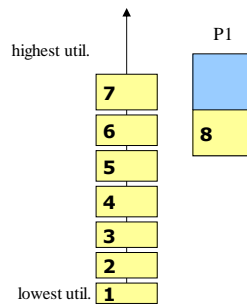
Lakshmanan's Algorithm [ECRTS'09]

- Pick up one processor, and assign as many tasks as possible



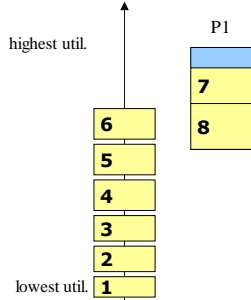
Lakshmanan's Algorithm [ECRTS'09]

- Pick up one processor, and assign as many tasks as possible



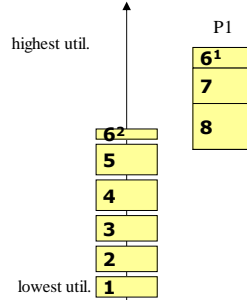
Lakshmanan's Algorithm [ECRTS'09]

- Pick up one processor, and assign as many tasks as possible



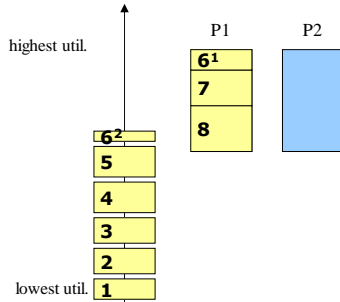
Lakshmanan's Algorithm [ECRTS'09]

- Pick up one processor, and assign as many tasks as possible



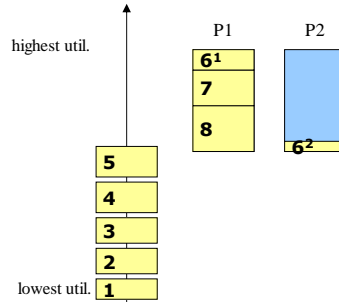
Lakshmanan's Algorithm [ECRTS'09]

- Pick up one processor, and assign as many tasks as possible



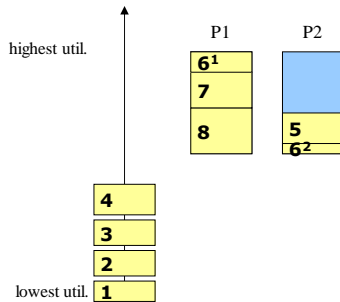
Lakshmanan's Algorithm [ECRTS'09]

- Pick up one processor, and assign as many tasks as possible



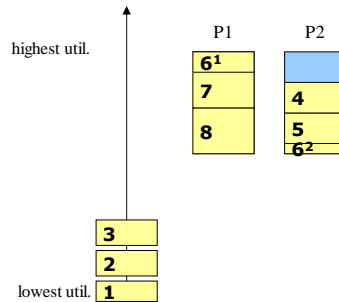
Lakshmanan's Algorithm [ECRTS'09]

- Pick up one processor, and assign as many tasks as possible



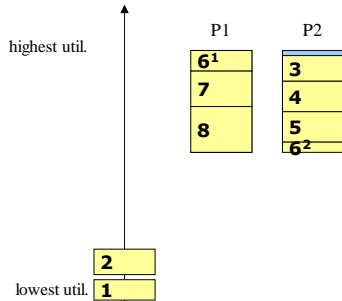
Lakshmanan's Algorithm [ECRTS'09]

- Pick up one processor, and assign as many tasks as possible



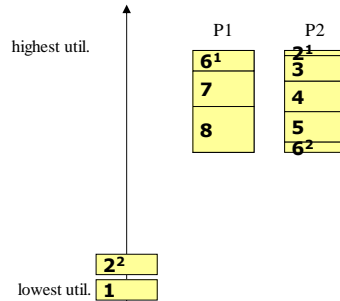
Lakshmanan's Algorithm [ECRTS'09]

- Pick up one processor, and assign as many tasks as possible



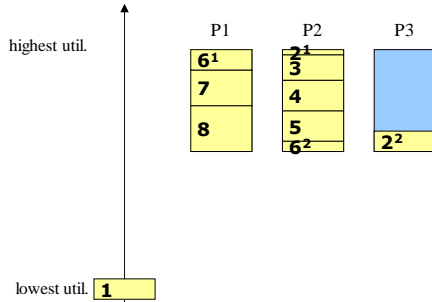
Lakshmanan's Algorithm [ECRTS'09]

- Pick up one processor, and assign as many tasks as possible



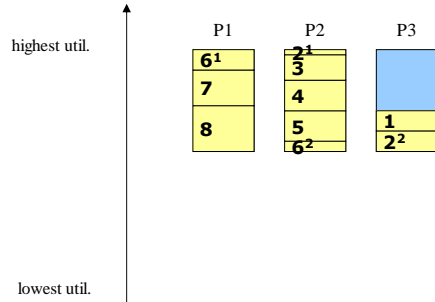
Lakshmanan's Algorithm [ECRTS'09]

- Pick up one processor, and assign as many tasks as possible



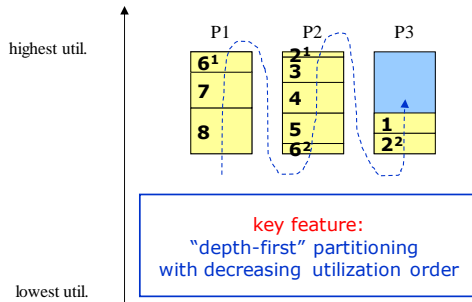
Lakshmanan's Algorithm [ECRTS'09]

- Pick up one processor, and assign as many tasks as possible



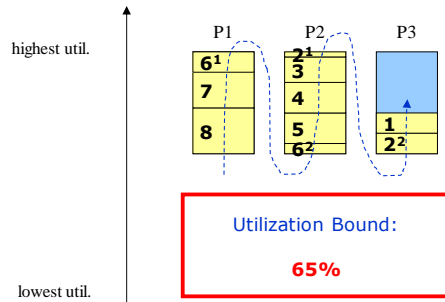
Lakshmanan's Algorithm [ECRTS'09]

- Pick up one processor, and assign as many tasks as possible



Lakshmanan's Algorithm [ECRTS'09]

- Pick up one processor, and assign as many tasks as possible

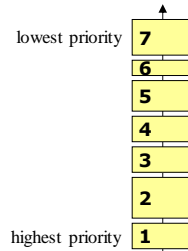


Our Algorithm [RTAS10]

"width-first" partitioning
with increasing priority order

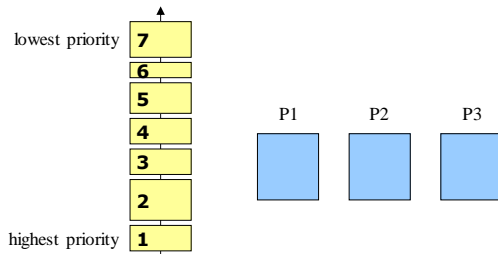
Our Algorithm

- Sort all tasks in **increasing priority order**



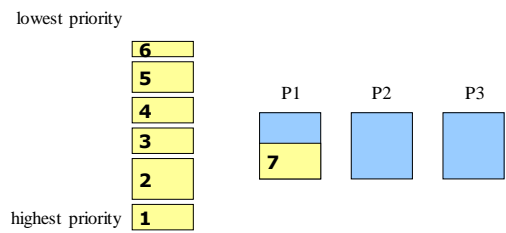
Our Algorithm

- Select the processor on which the assigned utilization is the **lowest**



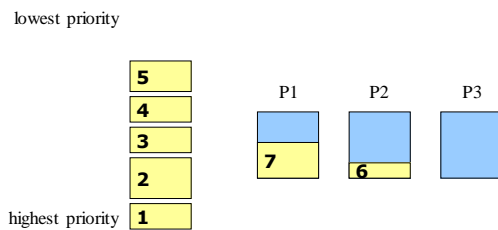
Our Algorithm

- Select the processor on which the assigned utilization is the **lowest**



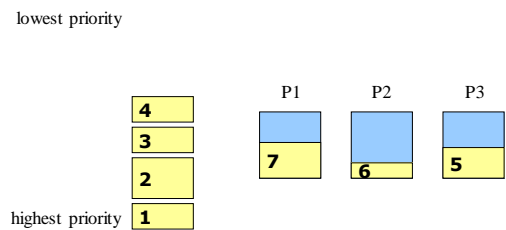
Our Algorithm

- Select the processor on which the assigned utilization is the **lowest**



Our Algorithm

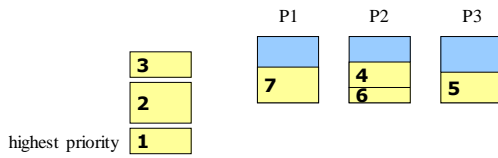
- Select the processor on which the assigned utilization is the **lowest**



Our Algorithm

- Select the processor on which the assigned utilization is the **lowest**

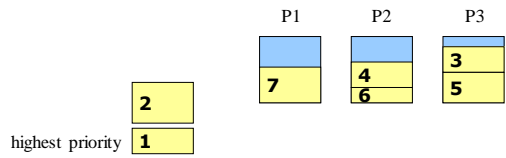
lowest priority



Our Algorithm

- Select the processor on which the assigned utilization is the **lowest**

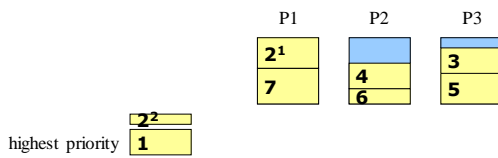
lowest priority



Our Algorithm

- Select the processor on which the assigned utilization is the **lowest**

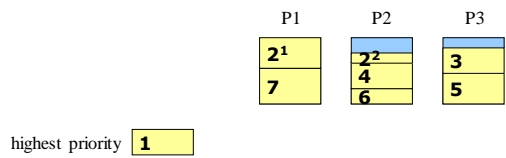
lowest priority



Our Algorithm

- Select the processor on which the assigned utilization is the **lowest**

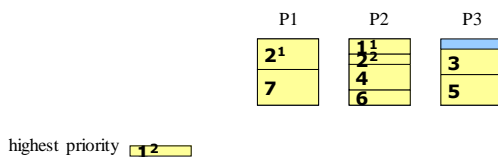
lowest priority



Our Algorithm

- Select the processor on which the assigned utilization is the **lowest**

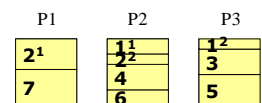
lowest priority



Our Algorithm

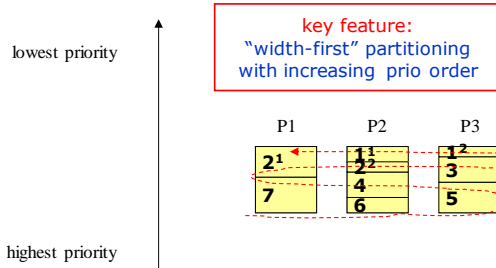
- Select the processor on which the assigned utilization is the **lowest**

highest priority 1²



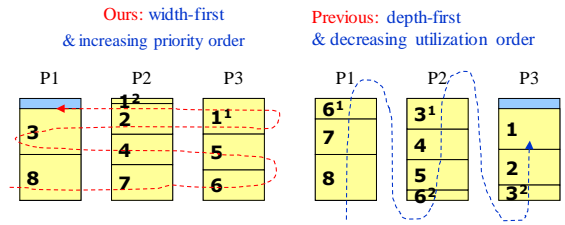
Our Algorithm

- Select the processor on which the assigned utilization is the **lowest**



Comparison

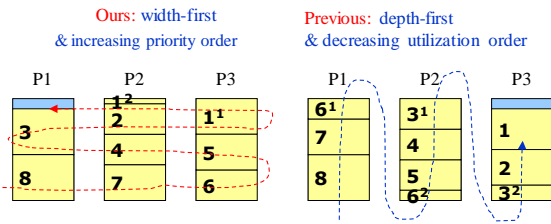
Why is our algorithm better?



Comparison

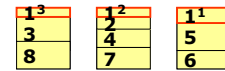
Why is our algorithm better?

By our algorithm split tasks generally have **higher priorities**



Split Task

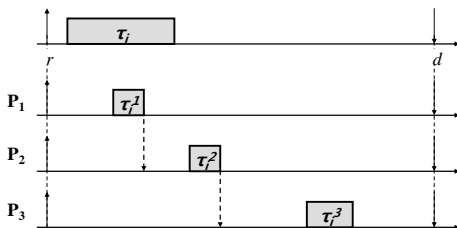
- Consider an extreme scenario:
 - suppose each subtask has the **highest** priority
 - schedulable anyway, we do not need to worry about their deadlines



- The difficult case is when the tail task is not on the top
 - the **key point is to ensure the tail task is schedulable**

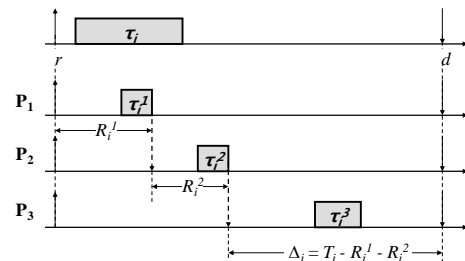
Split Task

- Subtasks should execute in the correct order



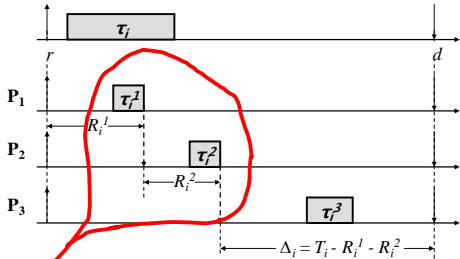
Split Task

- Subtasks get "shorter deadlines"



Split Task

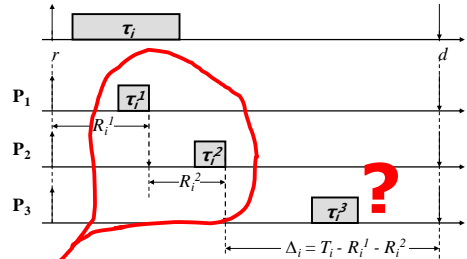
- Subtasks should execute in the correct order



These two are on the top: no problem with schedulability

Split Task

- Subtasks should execute in the correct order

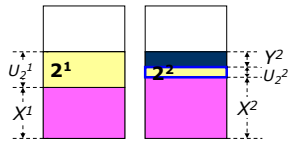


These two are on the top: no problem with schedulability

Why the tail task is schedulable?

The typical case: two CPUs and task 2 is split to two sub-tasks

As we always select the CPU with the lowest load assigned, we know



$$Y^2 + U_2^2 \leq U_2^1$$



$$Y^2 \leq U_2^1 - U_2^2$$

That is, the "blocking factor" for the tail task is bounded.

Theorem

For a task set in which each task τ_i satisfies

$$U_i \leq \frac{\Theta(N)}{1 + \Theta(N)}$$

we have

$$\frac{\sum C_i/T_i}{M} \leq N(2^{1/N} - 1)$$

\Rightarrow the task set is schedulable

$$\Theta(N) = N(2^{1/N} - 1) \quad N \rightarrow \infty, \quad \frac{\Theta(N)}{1 + \Theta(N)} \doteq 0.41$$

Theorem

For a task set in which each task τ_i satisfies

$$U_i \leq \frac{\Theta(N)}{1 + \Theta(N)} \quad \text{get rid of this constraint}$$

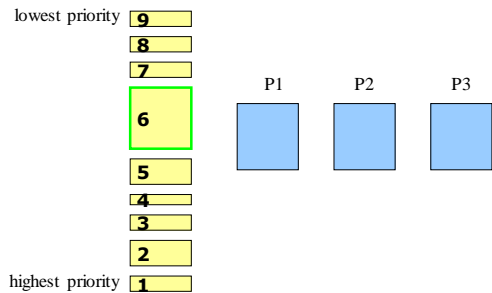
we have

$$\frac{\sum C_i/T_i}{M} \leq N(2^{1/N} - 1)$$

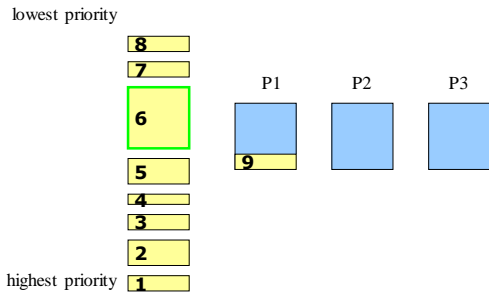
\Rightarrow the task set is schedulable

$$\Theta(N) = N(2^{1/N} - 1) \quad N \rightarrow \infty, \quad \frac{\Theta(N)}{1 + \Theta(N)} \doteq 0.41$$

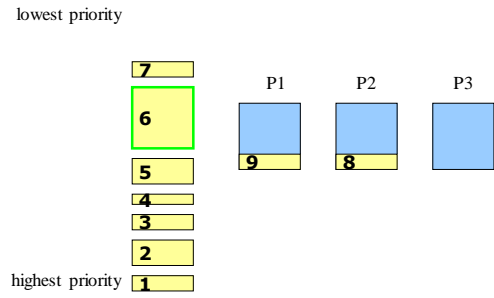
Problem of Heavy Tasks



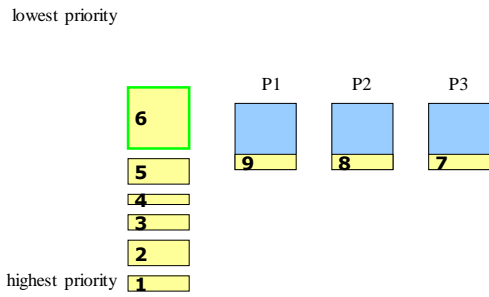
Problem of Heavy Tasks



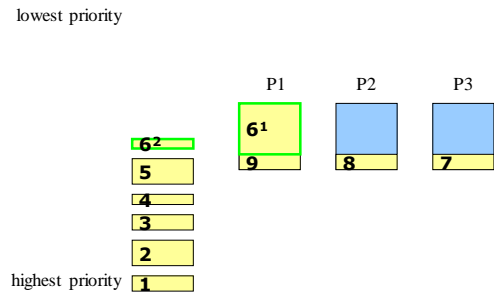
Problem of Heavy Tasks



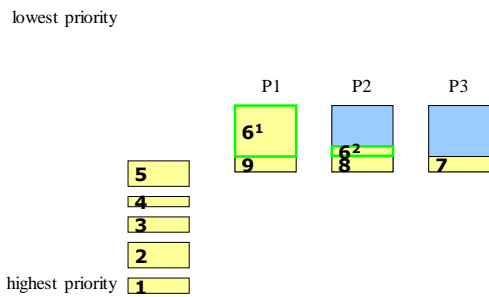
Problem of Heavy Tasks



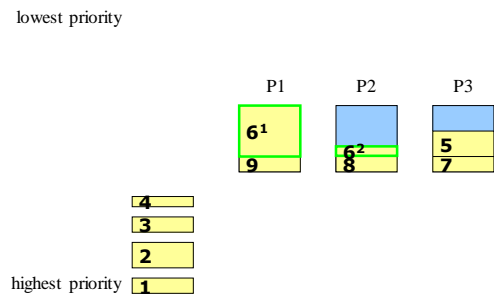
Problem of Heavy Tasks



Problem of Heavy Tasks

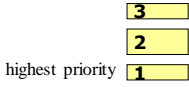
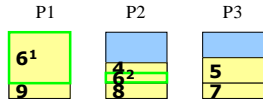


Problem of Heavy Tasks



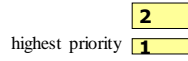
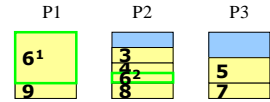
Problem of Heavy Tasks

lowest priority



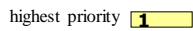
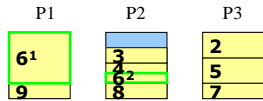
Problem of Heavy Tasks

lowest priority

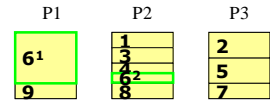


Problem of Heavy Tasks

lowest priority

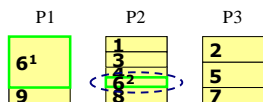


Problem of Heavy Tasks



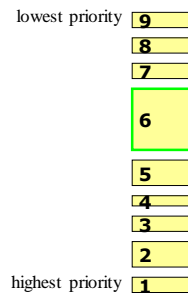
Problem of Heavy Tasks

the heavy tasks' tail task may have too low priority level



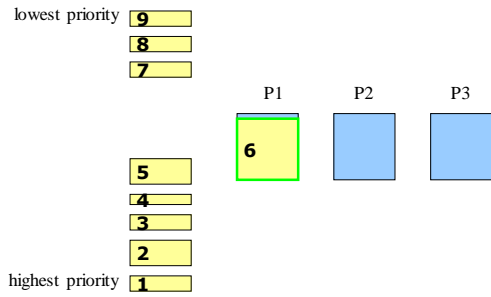
Solution for Heavy Tasks

Pre-assigning the heavy tasks (that may have low priorities)



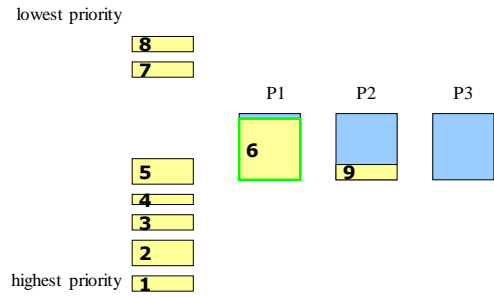
Solution for Heavy Tasks

- Pre-assigning the heavy tasks (that may have low priorities)



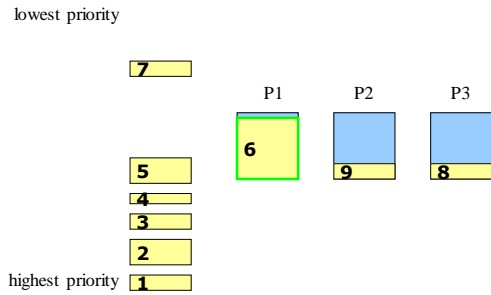
Solution for Heavy Tasks

- Pre-assigning the heavy tasks (that may have low priorities)



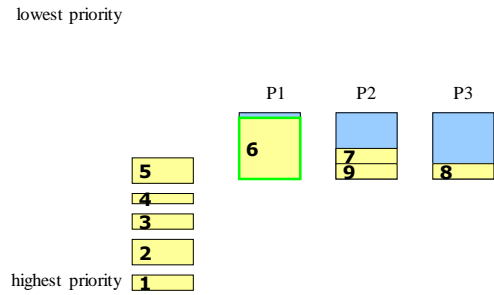
Solution for Heavy Tasks

- Pre-assigning the heavy tasks (that may have low priorities)



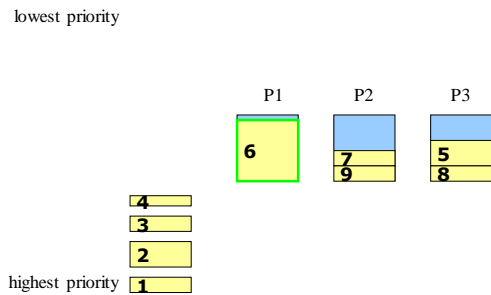
Solution for Heavy Tasks

- Pre-assigning the heavy tasks (that may have low priorities)



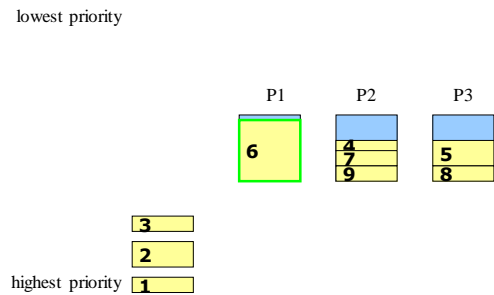
Solution for Heavy Tasks

- Pre-assigning the heavy tasks (that may have low priorities)



Solution for Heavy Tasks

- Pre-assigning the heavy tasks (that may have low priorities)



Solution for Heavy Tasks

- Pre-assigning the heavy tasks (that may have low priorities)
- lowest priority

P1	P2	P3
6	3	5
	4	8
	7	
	9	

highest priority 2
1

Solution for Heavy Tasks

- Pre-assigning the heavy tasks (that may have low priorities)
- lowest priority

P1	P2	P3
6	3	2
	4	5
	7	8
	9	

highest priority 1

Solution for Heavy Tasks

- Pre-assigning the heavy tasks (that may have low priorities)
- lowest priority

P1	P2	P3
6	1	2
	3	5
	4	8
	7	
	9	

highest priority 1 2

Solution for Heavy Tasks

- Pre-assigning the heavy tasks (that may have low priorities)

P1	P2	P3
1 2	1	2
6	3	5
	4	8
	7	
	9	

Solution for Heavy Tasks

- Pre-assigning the heavy tasks (that may have low priorities)

P1	P2	P3
1 2	1	2
6	3	5
	4	8
	7	
	9	

avoid to split heavy tasks (that may have low priorities)

Theorem

- By introducing the pre-assignment mechanism, we have

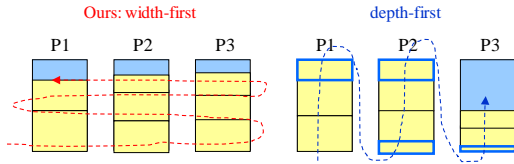
$$\frac{\sum C_i/T_i}{M} \leq N(2^{1/N} - 1)$$

⇒ the task set is schedulable

Liu and Layland's utilization bound for all task sets!

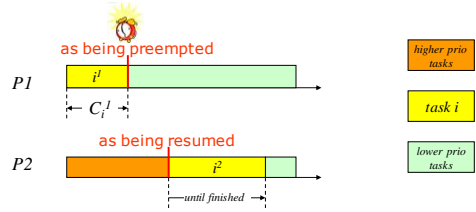
Overhead

- In both previous algorithms and ours
 - The number of task splitting is at most $M-1$
 - ❖ task splitting -> extra "migration/preemption"
 - Our algorithm on average has **less** task splitting

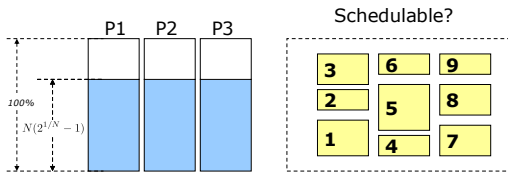


Implementation

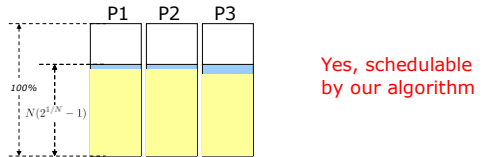
- Easy!
 - One timer for each split task
 - Implemented as "task migration"



Further Improvement

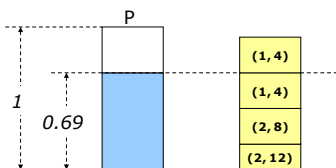


Using Liu and Layland's Utilization Bound



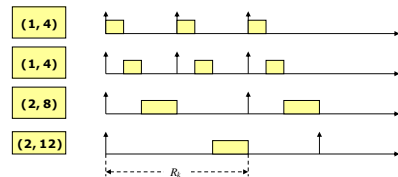
Utilization Bound is Pessimistic

- The Liu and Layland utilization bound is sufficient but not necessary
- many task sets are actually schedulable even if the total utilization is larger than the bound



Exact Analysis

- Exact Analysis: Response Time Analysis [Lehoczky_89]
 - pseudo-polynomial

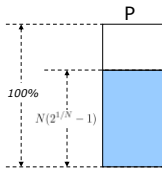
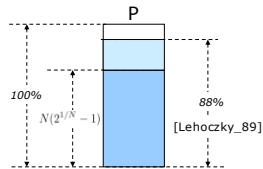


$$R_k = \sum_{T_i < T_k} \left\lceil \frac{R_k}{T_i} \right\rceil C_i + C_k$$

task k is schedulable iff $R_k \leq T_k$

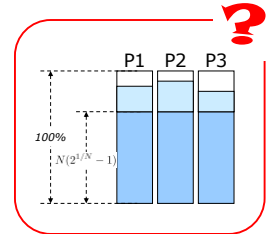
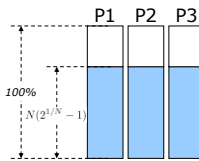
Utilization Bound v.s. Exact Analysis

- On single processors

Utilization bound Test
for RMSExact Analysis
for RMS

On Multiprocessors

- Can we do something similar on multiprocessors?

Utilization bound Test
the algorithm introduced above

Beyond Layland & Liu's Bound [RTSS 2010, rejected!]

- Our RTAS10 algorithm:
 - Increasing RMS priority order & worst-fit partitioning
 - Utilization test to determine the maximal load for each processor
 - The maximal load for each processor bounded by 69.3% $N(2^{1/N} - 1)$
- Improved algorithm:
 - Employ Response Time Analysis to determine the maximal workload on each processor
 - more flexible behavior (more difficult to prove ...)
 - Same utilization bound for the worst case, but
 - Much better average performance (by simulation)

I believe this is "the best algorithm" one can hope for "fixed-priority multiprocessor scheduling"

Conclusions

- The (multicore) Timing Problem is challenging
 - Difficult to guarantee Real-Time
 - and Difficult to analyze/predict
- Solutions: Partition & Isolation
 - Shared caches: coloring/partition
 - Memory bus/bandwidth: TDMA, ?
 - Processor cores: partition-based scheduling

Thanks!