

Multiprocessor Scheduling with Few Preemptions

Björn Andersson and Eduardo Tovar
IPP Hurray Research Group
Polytechnic Institute of Porto, Portugal
{bandersson,emt}@dei.isep.ipp.pt

Abstract

Consider the problem of scheduling a set of periodically arriving tasks on a multiprocessor with the goal of meeting deadlines. Processors are identical and have the same speed. Tasks can be preempted and they can migrate between processors. We propose an algorithm with a utilization bound of 66% and with few preemptions. It can trade a higher utilization bound for more preemptions and in doing so it has a utilization bound of 100%.

1. Introduction

Consider the problem of preemptive scheduling of n periodically arriving tasks on m identical processors. A task τ_i can arrive at any time when it arrives for the first time, but then it arrives periodically with a period T_i . Every time task τ_i arrives, it needs to execute C_i time units before it arrives again. A processor can execute at most one task at a time, and a task cannot execute on two or more processors simultaneously. The utilization is defined as $U_s = (1/m) \times \sum C_i / T_i$. The utilization bound UB_A of an algorithm A is the maximum number such that if $U_s \leq UB_A$ then all tasks meet their deadlines when scheduled by algorithm A .

The design space of preemptive multiprocessor scheduling algorithms can be categorized as *partitioned* or *global scheduling* [1, 2]. Global scheduling algorithms store in one queue (shared among all processors) all tasks that have arrived but have not finished their execution. At any time instant the m highest-priority tasks in that queue are selected for execution on the m processors using preemption and migration if necessary. In contrast, partitioned scheduling algorithms partition the set of tasks such that all tasks in a partition are assigned to the same processor. Tasks are not allowed to migrate from one processor to another processor, and hence the multiprocessor scheduling problem is transformed into m uniprocessor scheduling problems. This simplifies scheduling and schedulability analysis, since a large number of results available for uniprocessor scheduling can then be reused. Unfortunately, all partitioned

multiprocessor scheduling algorithms have a utilization bound of 50% or less [3]. Conversely, global scheduling can achieve a utilization bound of 100% by using a family of algorithms called *pfair scheduling* [4, 5]. Regrettably, this utilization bound comes at a price: all task parameters must be multiples of a time quantum, and in every time quantum a new task is selected for execution. As a result, the number of preemptions can be significantly high. We believe (as does Baker [6]) that it is desirable to achieve a higher (>50%) utilization bound without incurring the cost of an undesirable number of preemptions.

Therefore, in this paper we propose a multiprocessor scheduling algorithm with a utilization bound of 66%. It causes provably few preemptions: the number of preemptions divided by the number of jobs is at most 4. Of those algorithms in previous works that achieve the utilization bound (66%) of our algorithm, none of them has a finite number that bounds the number of preemptions divided by the number of jobs.

The remainder of this paper is organized as follows. Section 2 presents our new algorithm. Section 3 proves its utilization bound and Section 4 proves an upper bound on the number of preemptions per job. Section 5 offers discussion on previous work and conclusions.

2. The new algorithm

In order to understand the design of the new algorithm, we will first (in Section 2.1) consider partitioned scheduling of a specific task set example, with the purpose of stressing how partitioned scheduling may perform poorly. The reasoning on the example will provide the guiding principles on the design of the proposed new scheduling algorithm, which will then be formally presented in Sections 2.2 and 2.3.

2.1. Understanding the problem

Consider the following partitioned scheduling example: m processors and $n = m + 1$ tasks τ_i with $T_i = 1$ and $C_i = 0.5 + \varepsilon$. In partitioned scheduling, tasks cannot migrate; they are assigned to a processor and always execute there. Since $n > m$, there is one processor which is assigned two or more tasks. Therefore, the utilization of this processor will be at least $1 + 2\varepsilon$, and this is more than

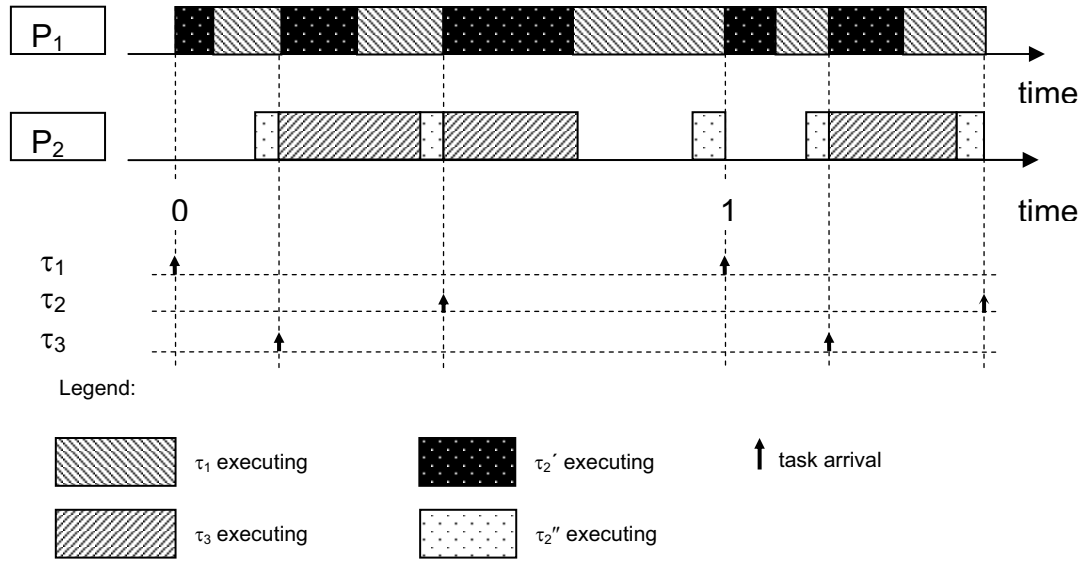


Figure 1: Schedule for the example task set using task splitting.

100%. Essentially, by choosing $m \rightarrow \infty$ and $\varepsilon \rightarrow 0$ the task set will have $U_s \rightarrow 0.5$ and, above all, a deadline is missed. Hence, the utilization bound of every partitioned scheduling algorithm is 50% or less.

This example stresses the fact that deadlines can be missed simply because a task could not be assigned to a processor, although there was plenty of idle time in the overall system. The idle time was spread out on different processors and could not be used. However, if in the same previous example a task is split into two sub-tasks and these sub-tasks are assigned to two different processors, then it is possible to assign tasks such that the utilization on every processor reaches 100%.

Note, however, that sub-tasks of the same task cannot execute simultaneously. A solution to this is a uniprocessor scheduling algorithm that is aware of tasks on other processors, so that a sub-task on one processor is not executed (even partially) at the same time as its corresponding sub-task on another processor.

Let us denote the two sub-tasks of a task τ_i as τ_i' and τ_i'' . Note that τ_i' and τ_i'' will have the same periods and will arrive at the same time. When a task (or sub-task) arrives on any processor, let t_0 denote the arrival time, and let t_1 denote the time of the next arrival of a job on any processor. Task τ_i' executes $(C_i' / T_i') \times (t_1 - t_0)$ time units without preemption and starts executing at time t_0 . Task τ_i'' executes $(C_i'' / T_i'') \times (t_1 - t_0)$ time units without preemption and finishes execution at time t_1 . Note that τ_i' and τ_i'' execute on different processors and their executions will not overlap in time since the original task τ_i had $C_i / T_i \leq 1$ and hence $C_i' / T_i' + C_i'' / T_i'' \leq 1$. We schedule these sub-tasks with the highest priority. The other tasks are scheduled according to EDF and have lower priority than

the sub-tasks. Applying this approach to the same task set example, would correspond to the schedule illustrated in Figure 1. There are three tasks $\{\tau_1, \tau_2, \tau_3\}$, with $T_i = 1$ and $C_i = 0.5 + \varepsilon$. These tasks are scheduled on two processors ($m = 2$): P_1 and P_2 . The task τ_2 has been split into two sub-tasks, τ_2' and τ_2'' , with $T_2' = T_2'' = 1$ and $C_2' = 0.5 - \varepsilon$, $C_2'' = 2\varepsilon$. Using the task splitting approach the task set is schedulable while it would not be using partitioned scheduling.

In the approach illustrated in Figure 1, task τ_i' is always executed at the beginning of the time interval between any two job arrivals while task τ_i'' is executed at the end of that time interval. However, by mirroring the schedule of the split tasks in half of the time intervals fewer preemptions would result. Figure 2 illustrates this for the same task set example.

2.2. Task assignment

We will now describe a general algorithm for assigning tasks to processors and scheduling tasks on a uniprocessor. The algorithm for assigning tasks to processors is given as pseudo-code in Algorithm 1 (Figure 3). The algorithm assigns tasks to processors such that on all processors the utilization does not exceed 100%. It has a parameter k which should be selected by the designer such that $1 \leq k \leq m$. The algorithm treats *heavy* and *light* tasks differently. A task τ_i is heavy if $C_i / T_i > SEP$, otherwise it is light. *SEP* means separator. The value of *SEP* is computed at line 4 in Algorithm 1, by calling a function described in Algorithm 3 (Figure 5). *SEP* will be used later on in Lemma 1 (Section 3), and the rationale for how to select it will be better perceived there. For the moment it is sufficient to realize that we select *SEP* as:

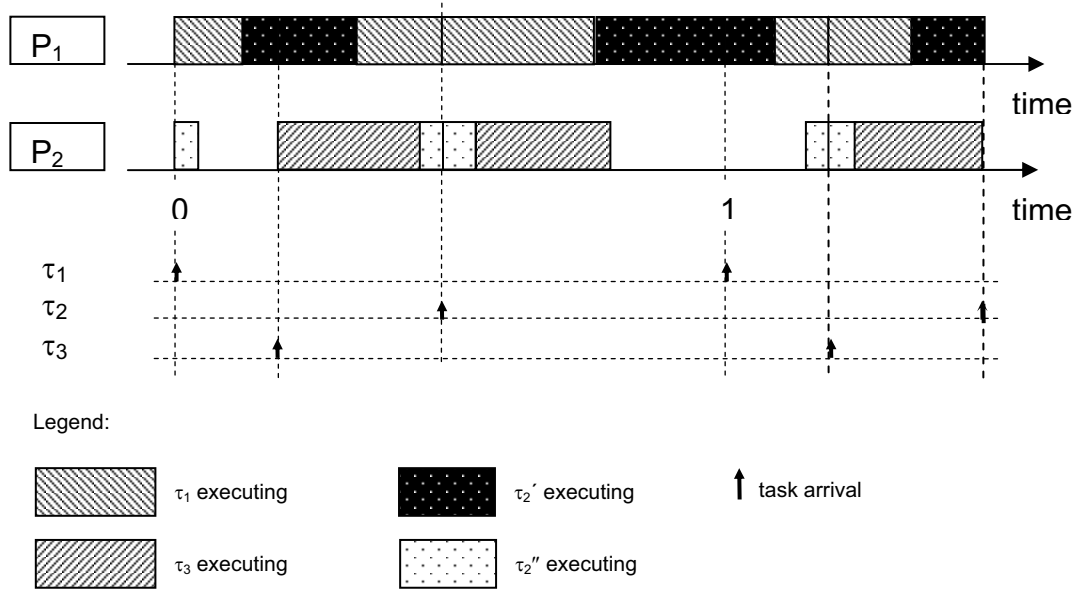


Figure 2: Using mirroring on the example in Figure 1 decreases the number of preemptions.

$$SEP = \begin{cases} \frac{k}{k+1} & k < m \\ 1 & k = m \end{cases} \quad (1)$$

First, the algorithm assigns heavy tasks to their own dedicated processors (at lines 10-13 in Algorithm 1). Doing so at this early stage in the algorithm improves performance, since the rest of it devoted to dealing with the light tasks. The main idea of the algorithm is that there is a current processor, with index p and tasks are considered one by one; index i denotes the current task. The task currently under consideration is attempted to be assigned to the current processor p . This is performed at line 18 in Algorithm 1. If the condition stated in that line (a schedulability test for EDF [7]) is true then the task can be assigned to processor p . If the condition is false, the task is split into two portions (at line 26) and assigned to the current processor p and processor $p + 1$ (lines 27-28). Then the processor with a higher index is considered (line 29).

As already pointed out, the algorithm assigns heavy tasks to some processors and light tasks to some other processors. L separates these tasks: heavy tasks are assigned to processors with indexes ranging from 1 up to L , while light tasks are assigned to processors with indexes ranging from $L + 1$ up to m . The light tasks are assigned to different groups of processors, and there are at most k processors in a group. Processors with indexes in the range $\{L + 1..L + k\}$ correspond to one group. Processors with indexes in the range $\{L + k + 1..L + 2k\}$ correspond to another group, and so forth. If a task is attempted to be assigned to the last processor in a group

and it fails, then it is not split; instead it is simply assigned to the next processor in a new group. Lines 23-24 take care of this situation in Algorithm 1. This ensures that tasks in one group do not interact with tasks in another group. This algorithm is actually similar to the *next-fit bin-packing* algorithm [2]. It differs however since task splitting is permitted in our approach.

2.3. Dispatching

We will now turn our attention to the problem of run-time dispatching. Our approach was already informally described in Section 2.1, and will now be formally described by Algorithm 2 (Figure 4). Algorithm 2 calls two subroutines, described in Algorithm 3 (Figure 5). All variables in Algorithm 2 and Algorithm 3 are global variables. The dispatching of heavy tasks is described on lines 1-7 in Algorithm 2. It entails a simple approach: whenever a task arrives, it executes on its assigned processor.

A light task is assigned to a group of processors, and whenever a task arrives in that group of processors, dispatchers are executed on all processors in that group. t_0 denotes the time when a task arrives, and t_1 denotes the time when any task in that group arrives next. This is computed on lines 20-21 within Algorithm 2. After that, Algorithm 2 calculates two time instants, $time_a$ and $time_b$, when tasks should be preempted. One of the split tasks is executed before $time_a$ (line 26) and the other split task is executed after $time_b$ (line 32). During the time span $[time_a, time_b)$ the non-split tasks are scheduled according to EDF. This is performed by the subroutine `schedule_tasks_with_EDF` described within Algorithm 3. If a non-split task finishes its execution during the time span

```

1.  for p in 1..m do
2.    U[p] := 0   τ[p] := {}
3.  end for
4.  SEP := calc_SEP
5.  Let τheavy denote the set of tasks with  $C_i/T_i > SEP$ 
6.  Let τlight denote the set of tasks with  $C_i/T_i \leq SEP$ 
7.  L := |τheavy|
8.  Order tasks such that τi with i in 1..L are all in τheavy and τi with i
   in L+1..n are all in τlight
9.  if |τheavy| ≤ m then
10.   for i in 1..L do
11.     p := i
12.     τ[p] := τ[p] ∪ {τi}      U[p] := U[p] + Ci/Ti
13.   end for
14.   if |τlight| > 0 then
15.     if L+1 ≤ m then
16.       p := L+1
17.       for i in L+1..n do
18.         if U[p]+Ci/Ti ≤ 1 then
19.           τ[p] := τ[p] ∪ {τi}   U[p] := U[p] + Ci/Ti
20.         else
21.           if p+1 ≤ m then
22.             if (p-L) mod k = 0 then
23.               p := p+1
24.               τ[p] := τ[p] ∪ {τi}   U[p] := U[p] + Ci/Ti
25.             else
26.               split task τi into τi' and τi" such that
27.                 Ti' = Ti
28.                 Ti" = Ti
29.                 Ci' = (1 - U[p]) * Ti'
30.                 Ci" = (Ci/Ti - Ci'/Ti') * Ti"
31.                 τ[p] := τ[p] ∪ {τi'}   U[p] := U[p] + Ci'/Ti'
32.                 τ[p+1] := τ[p+1] ∪ {τi"}   U[p+1] := U[p+1] + Ci"/Ti"
33.                 p := p+1
34.             end if
35.           else
36.             declare FAILURE
37.           end if
38.         end if
39.       end for
40.     else
41.       declare FAILURE
42.     end if
43.   else
44.     declare SUCCESS
45.   end if
46. end if

```

Figure 3: Algorithm 1 (for task assignment)

[time_a, time_b) then the next non-split task with the earliest deadline is selected. We denote the approach consisting of the use of Algorithms 1, 2 and 3 as *EKG approach*, as a short-hand notation for *EDF with task splitting and k processors in a group*.

3. Proving the utilization bound

The aim of this section is to prove that the utilization bound of EKG approach is *SEP* (as given by Equation (1), in Section 2.2). In order to do this, in Section 3.1, we prove that the task assignment scheme (Algorithm 1) declares success for all task sets with

utilization lower than or equal to the utilization bound. Lemma 2 states it and Lemma 1 is used to prove it. We also prove certain properties of the distribution of the execution of the task assignment in Lemmas 3-4. In Section 3.2 we show that if the properties assured by the task assignment scheme are true, then using the dispatcher (Algorithm 2) will enable all tasks to meet their deadlines. This proves the utilization bound and Theorem 1 states it.

```

1.  for processors with index p in 1..L:
2.    when the system starts do
3.      do nothing
4.    end do
5.    when a task arrives on processor p do
6.      execute that task
7.    end do
8.
9.  for processors with index p in L+1..m:
10.   when the system starts do
11.     mirrorflag[p] := false
12.     minindex[p] := L + ⌊(p-L)/k⌋ * k + 1
13.     maxindex[p] := L + ⌊(p-L)/k⌋ * k + k
14.     if maxindex[p] > m then
15.       maxindex[p] := m
16.     end if
17.     wait until all processors with index within range
18.       minindex[p]..maxindex[p] have reached this line
19.   end do
20.   when any task arrives on a processor with index in [minindex[p]..maxindex[p]] do
21.     t0 := current time
22.     t1 := next time a task arrives on a processor with index within range
23.       minindex[p]..maxindex[p]
24.     call calc_splitting_tasks // here we calculate timea_and_timeb
25.       // we also calculate firsttask and lasttask
26.     if firsttask[p] = NULL then
27.       idle processor p during [t0, timea[p]]
28.     else
29.       execute firsttask[p] on processor p during [t0, timea[p]]
30.     end if
31.     call schedule_tasks_with_EDF
32.     if lasttask[p] = NULL then
33.       idle processor p during [timeb[p], t1]
34.     else
35.       execute lasttask[p] on processor p during [timeb[p], t1]
36.     end if
37.     mirrorflag[p] := not (mirrorflag[p])
38.   end do

```

Figure 4: Algorithm 2 (run-time dispatching)

2.3. Task assignment

Lemma 1. If Algorithm 1 declares failure then the task set satisfies $U_s > SEP$.

Proof: One possibility is that the algorithm declared failure on line 43. If so, then all tasks that were assigned to processors had $C_i / T_i > SEP$. Since there are $m + 1$ or more tasks, it follows that:

$$\frac{1}{m} \sum_{i=1}^n \frac{C_i}{T_i} > \frac{1}{m} \times SEP \times (m+1) > SEP \quad (2)$$

Another possibility is that the algorithm declared failure on line 37. If this was the case then all m processors were assigned heavy tasks and then there was a light task in the task set that could not be assigned. This would imply the following:

$$\frac{1}{m} \sum_{i=1}^n \frac{C_i}{T_i} > \frac{1}{m} \times SEP \times m = SEP \quad (3)$$

Yet another possibility is that the algorithm declared failure on line 32. Then it must have been a task that was considered (call it τ_j) on line 18, and the

condition on that line was evaluated false. We observe that the utilization of the task set is no less than the sum of the utilization of all the processors when Algorithm 1 declared failure added to the utilization of τ_j . Hence:

$$\frac{1}{m} \sum_{i=1}^n \frac{C_i}{T_i} > \frac{1}{m} \times \left(\sum_{p=1}^m U[p] \right) + \frac{C_f}{T_f} \quad (4)$$

Inequality (4) can be rewritten to better express that the overall utilization is equal to the sum of the utilizations of each individual subset of processors:

$$\frac{1}{m} \sum_{i=1}^n \frac{C_i}{T_i} > \frac{1}{m} \times \left(\sum_{1 \leq p \leq L} U[p] + \sum_{(L+1 \leq p \leq m-1) \wedge ((p-L) \bmod k = 0)} U[p] + \sum_{(L+1 \leq p \leq m-1) \wedge ((p-L) \bmod k \neq 0)} U[p] + \sum_{p=m} U[p] + \frac{C_f}{T_f} \right) \quad (5)$$

```

1.  procedure calc_splitting_tasks is
2.  begin
3.    if mirrorflag[p] then
4.      if there is a task  $\tau_i'$  assigned
        to processor p then
5.        lasttask[p] :=  $\tau_i'$ 
6.        lastdur[p] :=  $(C_i'/T_i') * (t1-t0)$ 
7.      else
8.        lasttask[p] := NULL
9.        lastdur[p] := 0
10.     end
11.    if there is a task  $\tau_i''$  assigned
        to processor p then
12.      firsttask[p] :=  $\tau_i''$ 
13.      firstdur[p] :=  $(C_i''/T_i'') * (t1-t0)$ 
14.    else
15.      firsttask[p] := NULL
16.      firstdur[p] := 0
17.    end
18.  else
19.    if there is a task  $\tau_i'$  assigned
        to processor p then
20.      firsttask[p] :=  $\tau_i'$ 
21.      firstdur[p] :=  $(C_i'/T_i') * (t1-t0)$ 
22.    else
23.      firsttask[p] := NULL
24.      firstdur[p] := 0
25.    end
26.    if there is a task  $\tau_i''$  assigned
        to processor p then
27.      lasttask[p] :=  $\tau_i''$ 
28.      lastdur[p] :=  $(C_i''/T_i'') * (t1-t0)$ 
29.    else
30.      lasttask[p] := NULL
31.      lastdur[p] := 0
32.    end
33.  end if
34.  timea[p] :=  $t0 + firstdur[p]$ 
35.  timeb[p] :=  $t1 - lastdur[p]$ 
36. end
37.
38. procedure schedule_tasks_with_EDF is
39. begin
40.   let ready[p] denote the set of
        tasks  $\tau[p] \setminus \{firsttask[p], lasttask[p]\}$ 
        such that they have a job which
        has arrived at  $t0$  or earlier
        but the job has remaining
        execution at time  $t0$ 
41.   t[p] := timea[p]
42.   execute_entire_job[p] := true
43.   while execute_entire_job[p] do
44.     current_task[p] := dequeue the
        task with the least
        absolute deadline from
        ready[p]
45.     if current_task[p] = NULL then
46.       execute_entire_job[p] := false
47.     else
48.       remain[p] := the remaining execution
        time of the job of current_task
49.       if t[p] + remain[p] <= timeb[p] then
50.         execute_entire_job[p] := true
51.         execute the job of current task
52.         during [t[p], t[p] + remain[p]]
53.         t[p] := t[p] + remain[p]
54.       else
55.         execute_entire_job[p] := false
56.       end if
57.     end if
58.   end while
59.   if current_task[p] = NULL then
60.     idle processor p until timeb[p]
61.   else
62.     execute the job of current_task
        during [t[p], timeb[p]]
63.   end if
64. end
65.
66. function calc_SEP return real is
67. begin
68.   if  $k < m$  then
69.     return  $k / (k+1)$ 
70.   else
71.     return 1
72.   end if
73. end

```

Figure 5: Algorithm 3 (auxiliary procedures and functions)

Since all processors p with indexes $1 \leq p \leq L$ have $U[p] > SEP$ and $U[m] + C_f / T_f > 1$ (because the algorithm declared failure), inequality (5) can be re-written as follows:

$$\frac{1}{m} \sum_{i=1}^n \frac{C_i}{T_i} > \frac{1}{m} \times (L \times SEP + \sum_{(L+1 \leq p \leq m-1) \wedge ((p-L) \bmod k=0)} U[p] + \sum_{(L+1 \leq p \leq m-1) \wedge ((p-L) \bmod k \neq 0)} U[p] + 1) \quad (6)$$

Let us now consider the set of processors p with index $L+1 \leq p \leq m-1$. There are $(m-L-1)$ of these processors. For some of these processors, the condition at line 22 in Algorithm 1 resulted true. This happens when $(p-L) \bmod k = 0$. For the task

i , considered to be assigned to processor p , it must have been that $U[p] + C_i / T_i > 1$. Actually there are $\lfloor (m-L-1) / k \rfloor$ of such processors. Since task i is light then $U[p] > (1 - SEP)$. The other processors have $U[p] = 1$ when Algorithm 1 declared failure (the lines 26 and 27 guarantee that). Taking this reasoning into account, we can now state the following:

$$\frac{1}{m} \sum_{i=1}^n \frac{C_i}{T_i} > \frac{1}{m} \times (L \times SEP + \left\lfloor \frac{m-L-1}{k} \right\rfloor \times (1 - SEP) + m - L - 1 - \left\lfloor \frac{m-L-1}{k} \right\rfloor + 1) \quad (7)$$

We can obtain a lower bound on the right-hand side of (7). For details, see Appendix A in [8]. This gives us:

$$\frac{1}{m} \sum_{i=1}^n \frac{C_i}{T_i} > SEP \quad (8)$$

All cases where Algorithm 1 could have declared failure have been explored. If Algorithm 1 declares failure then inequalities (2), (3) and (8) enforce that $U_s > SEP$. This proves Lemma 1.

Lemma 2. If a task set satisfies $U_s \leq SEP$ and Algorithm 1 is used then Algorithm 1 declares success.

Proof: Follows from Lemma 1 and the fact that Algorithm 1 terminates on every input.

Lemma 3. If Algorithm 1 declares success then $\forall p: U[p] \leq 1$.

Proof: For those processors p with $p \leq L$, the claim that $U[p] \leq 1$ follows from the fact that $C_i / T_i \leq 1$ and the observation that Algorithm 1 lines 10-13 only assigns one task to each processor. For those processors p with $p \geq L + 1$, the claim $U[p] \leq 1$ follows from the actions taken by Algorithm 1 in line 18, line 24 and lines 26-28.

Lemma 4. If Algorithm 1 declares success then $\forall \tau_i$ that is split, $C_i' / T_i' + C_i'' / T_i'' \leq 1$.

Proof: Follows from line 26 in Algorithm 1 and the fact that before a task is split it satisfied the condition $C_i / T_i \leq 1$.

3.2. Dispatching

Lemma 5. Consider a time interval $[t_0, t_1]$ such that a task arrives at time t_0 and a task arrives at time t_1 but no tasks arrive during (t_0, t_1) . If $(\forall p: U[p] \leq 1)$ and (for all split tasks it holds that $C_i' / T_i' + C_i'' / T_i'' \leq 1$) and Algorithm 2 is used to schedule tasks on each processor in the group, then a task τ_i that was split by Algorithm 1 never executes on two or more processors simultaneously during the time span $[t_0, t_1]$.

Proof: Let p denote the processor to which τ_i' is assigned, and let $p + 1$ denote the processor to which τ_i'' is assigned. From Algorithm 2 it holds, for every processor p and every processor q , that $\text{mirrorflag}[p] = \text{mirrorflag}[q]$, since $\text{mirrorflag}[p]$ changes only when a task arrives. We will consider two cases.

Case 1. $\text{mirrorflag}[p] = \text{false}$ during the time span $[t_0, t_1]$.

From Algorithm 2, τ_i' will execute during $[t_0, \text{timea})$ on processor p and τ_i'' will execute during $[\text{timeb}, t_1]$ on processor $p + 1$. An assumption associated to the lemma is that $C_i' / T_i' + C_i'' / T_i'' \leq 1$. Hence:

$$0 \leq (t_1 - t_0) \times (1 - (C_i' / T_i' + C_i'' / T_i'')) \quad (9)$$

which can be re-written as follows:

$$0 \leq t_1 - (t_1 - t_0) \times C_i' / T_i' - (t_0 + (t_1 - t_0) \times C_i'' / T_i'') \quad (10)$$

From Algorithm 3 (lines 34-35), the two terms in (10) correspond to timea and timeb , and therefore it results that (10) can be re-written as follows:

$$0 \leq \text{timeb} - \text{timea} \quad (11)$$

Clearly, from (11), $\text{timea} \leq \text{timeb}$ and this assures that τ_i' and τ_i'' do not execute simultaneously during the time span $[t_0, t_1]$.

Case 2. $\text{mirrorflag}[p] = \text{true}$ the time span $[t_0, t_1]$

The reasoning is similar to Case 1.

Thus, regardless of which one of the two cases is true, it holds that a split task does not execute simultaneously on two or more processors.

Lemma 6. If $(\forall p: U[p] \leq 1)$ and (for all split tasks it holds that $C_i' / T_i' + C_i'' / T_i'' \leq 1$) and (Algorithm 2 is used to schedule tasks on each processor in the group) then a task never executes on two or more processors simultaneously.

Proof: We will prove this lemma by proving that an arbitrary task τ_i never executes on two or more processors simultaneously.

Case 1. The task τ_i was not split by Algorithm 1.

Case 1a). τ_i is assigned to a processor p with $p \leq L$. This task is assigned to a processor and does not execute on any other processor. Hence τ_i never executes on two or more processors simultaneously.

Case 1b). τ_i is assigned to a processor p with $L + 1 \leq p$. This case is similar to Case 1a); but on a processor p that satisfies $L + 1 \leq p$ there may be other tasks executing. Nevertheless, the reasoning in Case 1a) can be used in this case too, and thus supporting that a task τ_i never executes on two or more processors simultaneously.

Case 2. The task τ_i was split by Algorithm 1.

Consider a time interval $[t_0, t_1]$ such that a task arrives at time t_0 and a task arrives at time t_1 but no tasks arrive during (t_0, t_1) . From Lemma 5 it results that τ_i' and τ_i'' never execute simultaneously during $[t_0, t_1]$. This argument can be repeated for every time interval between any two consecutive task arrivals, and hence it results that τ_i never executes simultaneously on two or more processors.

Therefore, regardless of which one of the cases is true, the lemma holds.

Lemma 7. Consider a time interval $[t_0, t_1]$ such that a task arrives at time t_0 and a task arrives at time t_1 but no tasks arrive during (t_0, t_1) . If $(\forall p: \cup[p] \leq 1)$ and (for all split tasks it holds that $C_i' / T_i' + C_i'' / T_i'' \leq 1$) and (Algorithm 2 is used to schedule tasks on each processor in the group) then a task τ_i that was split by Algorithm 1 executes $(C_i / T_i) \times (t_1 - t_0)$ time units during $[t_0, t_1]$.

Proof: Let p denote the processor to which τ_i' is assigned, and let $p + 1$ denote the processor to which τ_i'' is assigned. From Algorithm 2, and for every processor p and every processor q , it holds that $\text{mirrorflag}[p] = \text{mirrorflag}[q]$ since $\text{mirrorflag}[p]$ changes only when a task arrives.

We will consider two cases.

Case 1. $\text{mirrorflag}[p] = \text{false}$ during the time span $[t_0, t_1]$.

From Algorithm 2, it results that τ_i' executes during $[t_0, \text{timea})$ on processor p and τ_i'' executes during $[\text{timeb}, t_1)$ on processor $p + 1$. Hence, the amount that τ_i executes during $[t_0, t_1]$ is given by:

$$(t_1 - \text{timeb}) + (\text{timea} - t_0) \quad (12)$$

As reasoned previously in the proof of Lemma 5, Algorithm 2 states that:

$$\begin{cases} \text{timea} = t_0 + (t_1 - t_0) \times C_i' / T_i' \\ \text{timeb} = t_1 - (t_1 - t_0) \times C_i'' / T_i'' \end{cases} \quad (13)$$

and hence, (12) can be re-written as follows:

$$(t_1 - t_0) \times C_i' / T_i' + (t_1 - t_0) \times C_i'' / T_i'' \quad (14)$$

Since from Algorithm 1 (line 26) a split task τ_i was split such that $C_i' / T_i' + C_i'' / T_i'' = C_i / T_i$ and $T_i' = T_i'' = T_i$, then the amount of execution performed by τ_i during the time interval $[t_0, t_1]$ is given by:

$$(t_1 - t_0) \times C_i / T_i \quad (15)$$

Case 2. $\text{mirrorflag}[p] = \text{true}$ during the time span $[t_0, t_1]$

The reasoning is similar to Case 1.

Thus, and regardless of which case occurs, the statement of the lemma is obtained.

Lemma 8. If $(\forall p: \cup[p] \leq 1)$ and (for all split tasks it holds that $C_i' / T_i' + C_i'' / T_i'' \leq 1$) and (Algorithm 2 is used

to schedule tasks on each processor in the group) then deadlines are met for all tasks τ_i such that $(\tau_i$ is assigned to a processor p with $L + 1 \leq p) \wedge (\tau_i$ is not split).

Proof: We will make definitions and algebraic manipulations in the first part of the proof and use them to prove the lemma by contradiction in a second part of the proof. 2.

1. Let $\tau^{\text{not_split}}[p]$ denote the set of tasks such that each task in $\tau^{\text{not_split}}[p]$ is $(\tau_i$ is assigned a processor p with $L + 1 \leq p) \wedge (\tau_i$ is not split). Let τ_a' denote the prime-task assigned to processor p and let τ_b'' denote the bis-task assigned to processor p . Note that τ_a' and τ_b'' do not belong to the same original task. From Algorithm 1 it results that:

$$\left(\sum_{i \in \tau^{\text{not_split}}[p]} \frac{C_i}{T_i} \right) + \frac{C_a'}{T_a'} + \frac{C_b''}{T_b''} \leq 1 \quad (16)$$

which can be re-written as follows:

$$\left(\sum_{i \in \tau^{\text{not_split}}[p]} \frac{C_i}{T_i} \right) \leq 1 - \left(\frac{C_a'}{T_a'} + \frac{C_b''}{T_b''} \right) \quad (17)$$

Both sides of inequality (17) can be multiplied by S , where S is any positive real number. In doing it and applying an algebraic re-writing, the following inequality holds:

$$\begin{aligned} \forall S > 0: \quad & \sum_{i \in \tau^{\text{not_split}}[p]} \left(\left\lfloor \frac{S - T_i}{T_i} \right\rfloor + 1 \right) \times C_i \\ & \leq S \times \left(1 - \left(\frac{C_a'}{T_a'} + \frac{C_b''}{T_b''} \right) \right) \end{aligned} \quad (18)$$

2. Let us now consider the task set $\tau^{\text{not_split}}[p]$, and prove that all tasks in $\tau^{\text{not_split}}[p]$ meet their deadlines. We prove it using contradiction. Suppose that a task in $\tau^{\text{not_split}}[p]$ missed a deadline. Let us study the busy period [9] before the first deadline miss. Our busy period is defined slightly different from the definition in [9]; our busy period starts when a task arrives and ends when the first deadline miss occurs. During this time period, it is permitted that processor p is idle. However it must be that $\text{ready}[p]$ (see line 40 in Algorithm 3) is non-empty. Just before the beginning of the busy period, $\text{ready}[p]$ (line 40 in Algorithm 3) is empty. Since the tasks in $\tau^{\text{not_split}}[p]$ are scheduled with preemptive EDF (pseudo-code in Algorithm 3 – Figure 5) it results that for a deadline miss to occur, it must have been that there is a time interval Q such that Q is a

subset of our busy period) $\wedge (Q$ starts when a task arrives and ends when the deadline is missed) \wedge (the supply of processing time [9] for the tasks in $\tau^{not_split}[p]$ during Q was less than the demand of processing time from tasks in $\tau^{not_split}[p]$ during Q). Let LQ denote the length of Q . It is known that the demand is calculated as follows:

$$\sum_{i \in \tau^{not_split}[p]} \left(\left\lfloor \frac{LQ - T_i}{T_i} \right\rfloor + 1 \right) \times C_i \quad (19)$$

Let us consider a time interval $[t_0, t_1)$ such that at time t_0 , a task arrives and at time t_1 a task arrives but during (t_0, t_1) no tasks arrive. Let us consider a task τ_a' and a task τ_b' which are split tasks and are assigned to the same processor p . Note that they do not belong to the same original task. It holds that the task τ_a' executes for $(C_a' / T_a') \times (t_1 - t_0)$ time units and the task τ_b' executes for $(C_b'' / T_b'') \times (t_1 - t_0)$ time units. This holds for every such time $[t_0, t_1)$. Hence, the supply of processing time for the tasks in $\tau^{not_split}[p]$ in the time interval $[t_0, t_1)$ is given by:

$$(t_1 - t_0) \times \left(1 - \left(\frac{C_a'}{T_a'} + \frac{C_b''}{T_b''} \right) \right) \quad (20)$$

This is true regardless of the value of $mirror[p]$. Observe that the time interval Q can be subdivided into intervals where inequality (20) can be applied. The repeated application of (20) to all those intervals yields that the supply of processing time for the task in $\tau^{not_split}[p]$ during Q is:

$$LQ \times \left(1 - \left(\frac{C_a'}{T_a'} + \frac{C_b''}{T_b''} \right) \right) \quad (21)$$

It is known from previous research [9] that since a deadline is missed the demand exceeded the supply. Hence the following inequality will hold:

$$\sum_{i \in \tau^{not_split}[p]} \left(\left\lfloor \frac{LQ - T_i}{T_i} \right\rfloor + 1 \right) \times C_i > LQ \times \left(1 - \left(\frac{C_a'}{T_a'} + \frac{C_b''}{T_b''} \right) \right) \quad (22)$$

But this contradicts (18). Hence, all deadlines of tasks in $\tau^{not_split}[p]$ are met.

Lemma 9. If $(\forall p: \cup[p] \leq 1)$ and (for all split tasks it holds that $C_i' / T_i' + C_i'' / T_i'' \leq 1$) and (Algorithm 2 is used to schedule tasks on each processor in the group) then all deadlines are met.

Proof: We will prove this lemma by proving that an arbitrary task τ_i meets its deadline.

Case 1. The task τ_i was not split in Algorithm 1.

Case 1a). τ_i is assigned to a processor p with $p \leq L$.

There are no other tasks on this processor. Since $C_i \leq T_i$, the task meets its deadlines.

Case 1b). τ_i is assigned to a processor p with $L + 1 \leq p$.

It results from Lemma 8 that τ_i meets its deadlines.

Case 2. The task τ_i was split in Algorithm 1.

Consider a time interval $[t_0, t_1)$ such that a task arrives at time t_0 and a task arrives at time t_1 but no tasks arrive during (t_0, t_1) . From Lemma 7 it results that τ_i executes $(C_i / T_i) \times (t_1 - t_0)$ time units during $[t_0, t_1)$. This argument can be repeated for every time interval between two task arrivals, and hence it results that τ_i meets all its deadlines.

Therefore, regardless of which one of the cases is true, the lemma holds.

We will finalise Section 3 by stating the following theorem.

Theorem 1. If a task set satisfies the condition $U_s \leq SEP$ and EKG is used, then all deadlines are met and a task never executes on two or more processors simultaneously.

Proof: It follows from Lemma 6 and Lemma 9.

4. Proving a bound on the number of preemptions

In this section we will state (Theorem 2) a bound on the number of preemptions divided by the number of jobs for the entire task set. Its proof depends on Lemma 10, which studies the special case of tasks in one group of processors. We let lcm denote the least common multiple of periods of the tasks. Since Lemma 10 and Theorem 2 makes claims about the number of preemptions per job, it is necessary to define whether a split task is one task or two tasks; analogously for jobs. In this section, in Lemma 10 and Theorem 2, we say that if a task is split then it is counted as one task. The jobs from these tasks are split. We count such a split job as one job. We think this is the right way of counting because the splitting of tasks and jobs is only used internally in the scheduling algorithm; the application does not know that a task or a job is split.

Lemma 10. Consider the case when all tasks arrive at time 0 and they are scheduled using EKG. Consider the processors in one of the groups consisting of k processors. We claim that

the number of preemptions during $[0, lcm)$ divided by the number of jobs that executed during $[0, lcm)$ is at most $2k$.

Proof: See Appendix B in [8]. \square

Theorem 2. Consider the case when all tasks arrive at time 0 and they are scheduled using EKG. We claim that the number of preemptions during $[0, lcm)$ divided by the number of jobs that executed during $[0, lcm)$ is at most $2k$.

Proof: From Lemma 10 we know that for the groups with k processors, the number of preemptions divided by the number of jobs is at most $2 \times k$. The last group may consist of less than k processors; nevertheless, it holds that the number of preemptions divided by the number of jobs is at most $2 \times k$. The processors with index $p \leq L$ have only one task per processor so no preemptions can occur there. Taking these three facts together imply that the number of preemptions on all processors during $[0, lcm)$ divided by the number of jobs during $[0, lcm)$ in the entire task set is no greater than $2 \times k$.

5. Conclusions

We have presented an algorithm to schedule tasks to meet deadlines on a multiprocessor. By selecting $k = 2$, Theorem 1 tells us that the utilization bound of the algorithm is 66% and Theorem 2 tells us that it causes at most 4 preemptions per job on average per hyperperiod. A higher value of k , gives a higher utilization bound at the expense of more preemptions. By selecting $k = m$ we obtain that the utilization bound is 100%.

Pfair scheduling algorithms [4, 5] and the algorithm BF [10] have a utilization bound of 100%. Unfortunately neither of them have proven bounds on the number of preemptions per job. Several algorithms were presented by Khemka and Shyamasundar for multiprocessor scheduling with a high utilization bound and few preemptions [11]. They offered a bound on the number of preemptions as a function of the greatest common divisor of the periods of all tasks in the task set. Unfortunately, for some task sets the bound on the number of preemptions divided by the number of jobs approaches infinity. This happens for task sets with periods selected as $T_i = (q + i)$:th prime number and where q approaches infinity. As already pointed out, with our algorithm this cannot happen.

Acknowledgements

This work was partially funded by Fundação para Ciência e Tecnologia (FCT).

References

- [1] J. Leung and J. Whitehead, "On the Complexity of Fixed-priority Scheduling of Periodic Real-Time Tasks," *Performance Evaluation, Elsevier Science*, vol. 22, pp. 237-250, 1982.

- [2] S. K. Dhall and C. L. Liu, "On a real-time scheduling problem," *Operations Research*, vol. 26, pp. 127-140, 1978.
- [3] D.-I. Oh and T. P. Baker, "Utilization Bounds for N-Processor Rate Monotone Scheduling with Static Processor Assignment," *Real Time Systems Journal*, vol. 15, pp. 183-192, 1998.
- [4] S. K. Baruah, N. K. Cohen, C. G. Plaxton, and D. A. Varvel, "Proportionate Progress: A Notion of Fairness in Resource Allocation," *Algorithmica*, vol. 15, pp. 600-625, 1996.
- [5] J. Anderson and A. Srinivasan, "Mixed Pfair/ERfair Scheduling of Asynchronous Periodic Tasks," *Journal of Computer and System Sciences*, vol. 68, pp. 157-204, 2004.
- [6] T. P. Baker, "Comparison of Empirical Success Rates of Global vs. Partitioned Fixed-Priority EDF Scheduling for Hard Real Time," Department of Computer Science, Florida State University, Tallahassee, FL 32306 July 2005.
- [7] C. L. Liu and J. W. Layland, "Scheduling Algorithms for Multiprogramming in a Hard-Real-Time Environment," *Journal of the ACM (JACM)*, vol. 20, pp. 46-61, 1973.
- [8] B. Andersson and E. Tovar, "Multiprocessor Scheduling with Few Preemptions," IPP Hurray Research Group, Polytechnic Institute of Porto, Portugal HURRAY-TR-060501, Available at <http://www.hurray.isep.ipp.pt/privfiles/tr-hurray-060501.pdf>, May 2006.
- [9] S. K. Baruah, R. Howell, and L. Rosier, "Algorithms and complexity concerning the preemptive scheduling of periodic, real-time tasks on one processor," *Real-Time Systems*, pp. 301-324, 1990.
- [10] D. Zhu, D. Mossé, and R. Melhem, "Multiple-Resource Periodic Scheduling Problem: how much fairness is necessary?," presented at 24th IEEE International Real-Time Systems Symposium, 2003.
- [11] A. Khemka and R. K. Shyamasundar, "Multiprocessor Scheduling of Periodic Tasks in a Hard Real-Time Environment," presented at International Parallel Processing Symposium, 1992.