

Document downloaded from:

<http://hdl.handle.net/10251/38465>

This paper must be cited as:

March Cabrelles, J.L.; Sahuquillo Borrás, J.; Petit Martí, S.V.; Hassan Mohamed, H.; Duato Marín, J.F. (2013). Power-aware scheduling with effective task migration for real-time multicore embedded systems. *Concurrency and Computation: Practice and Experience*. 25(14):1987-2001. doi:10.1002/cpe.2899.



The final publication is available at

<http://onlinelibrary.wiley.com/doi/10.1002/cpe.2899/pdf>

Copyright Wiley

Power-Aware Scheduling with Effective Task Migration for Real-Time Multicore Embedded Systems

J. L. March^{1*}, J. Sahuquillo¹, S. Petit¹, H. Hassan¹ and J. Duato¹

¹*Department of Computer Engineering (DISCA), Universitat Politècnica de València, Spain*

SUMMARY

A major design issue in embedded systems is reducing the power consumption since batteries have a limited energy budget. For this purpose, several techniques such as Dynamic Voltage Scaling (DVS) or task migration are being used. DVS allows reducing power by selecting the optimal voltage supply, while task migration achieves this effect by balancing the workload among cores.

This paper focuses on power-aware scheduling allowing task migration to reduce energy consumption in multicore embedded systems implementing DVS capabilities. To address energy savings, the devised schedulers follow two main rules: migrations are allowed at specific points of time and only one task is allowed to migrate each time.

Two algorithms have been proposed working under real-time constraints. The simpler algorithm, namely, Single Option Migration (SOM) only checks just one target core before performing a migration. In contrast, the Multiple Option Migration (MOM) searches the optimal target core.

In general, the MOM algorithm achieves better energy savings than the SOM algorithm, although differences are wider for a reduced number of cores and frequency/voltage levels. Moreover, the MOM algorithm reduces energy consumption as much as 40% over the typical WF algorithm.

Copyright © 2011 John Wiley & Sons, Ltd.

Received ...

KEY WORDS: Dynamic Partitioning; Task Migration; Embedded Systems; Power-Aware; Multicore; Real-Time Systems

1. INTRODUCTION

Embedded systems are an important segment of the microprocessor market as they are becoming ubiquitous in our life. Systems like PDAs, smartphones, and automotive devices, provide an increasing number of functionalities such as navigation, multimedia or gaming, so that computational power is becoming more important every day. However, increasing computational power impacts on battery lifetime, so a major design concern is power optimization and management [1, 2].

*Correspondence to: Department of Computer Engineering (DISCA), Universitat Politècnica de València, Camino de Vera s/n, 46022, Valencia, Spain. E-mail: jomarcab@gap.upv.es

To deal with both computational and power management requirements, many systems use multicore processors. These processors allow a more efficient power management than complex monolithic processors for a given performance level. Moreover, many manufacturers (e.g., Intel, IBM, Sun, etc.) deliver processors providing multithreading capabilities, that is, they provide support to run several threads simultaneously. Some examples of current multithreaded processors are Intel Montecito [3] and IBM Power 5 [4]. Also, leading manufacturers of the embedded sector, like ARM, plan to include multithreading technology in next-generation processors [5].

A power management technique that is being implemented in most current microprocessors is Dynamic Voltage Scaling (DVS) [6]. This technique allows the system to improve its energy consumption by reducing the frequency when the processor has a low activity level (e.g., a mobile phone that is not being actively used). In a multicore system, the DVS regulator can be shared by several cores, also referred to as *global*, which means that leakage power consumption is mostly the same in all the cores. On the contrary, some systems have a local or *private* DVS regulator for each individual core. In the former case, all cores are forced to work at the same speed but less regulators are required so it is a cheaper solution. The latter case could enable more energy savings since each core frequency can be properly tuned to its applications requirements but it is more expensive [7].

On the other hand, energy consumption in systems with a global DVS regulator can be further improved by properly balancing the workload [1, 8]. To this end, a *partitioner* module is in charge of distributing the workload (i.e., the set of tasks) according to a given algorithm, such as the Worst Fit (WF) or First Fit [9], that selects the target core to run each task. Unfortunately, the nature of some workload mixes prevents the partitioner from achieving a good balance. To deal with this drawback some systems allow tasks to migrate and move their execution from one core to another, which results in energy saving improvements.

In this paper, two algorithms allowing task migration to reduce energy consumption in multicore embedded systems with real-time constraints implementing DVS capabilities, are proposed. The simpler algorithm, namely, Single Option Migration (SOM) only checks just one target core before performing a migration. In contrast, the Multiple Option Migration (MOM) searches the optimal target core. To address energy savings, the devised schedulers follow two main rules: (i) migrations are allowed at specific points of time because analyzing all the possible task migrations may result in a prohibitive overhead, (ii) and only one task is allowed to migrate each time. This work focuses on multicore processors where the scheduler includes a partitioner module to distribute tasks among cores. This partitioner is in charge of readjusting possible workload imbalances at run-time that may occur at arrivals or exits of tasks by applying task migration. To keep overhead low and studying the impact of the point of time when the algorithm is applied, three variants of the SOM algorithm have been devised, depending on the point of time the scheduler is applied: at task arrival to the system (SOM_{in}), when a task leaves the system (SOM_{out}), and in both cases (SOM_{in-out}).

Because of energy constraints, embedded systems are still limited to a lower number of cores than their high-performance counterparts. Therefore, energy evaluation results focus on a realistic number of cores: two, three and four cores. Some examples are the bi-core Intel Atom [10], the tri-core Marvell ARMADATM 628 [11] or the quad-core ARM 11 MPCore [12]. On the other hand, this work assumes a relatively wide number of frequency/voltage levels (up to eight) in order to approach the results to real systems.

Experimental results show that applying the algorithm at tasks' exits can achieve better energy savings than applying only at tasks' arrivals, but the highest benefits are obtained when the algorithm is applied in both cases. The MOM algorithm achieves better energy savings than the SOM algorithm. Differences are wider for a reduced number of cores and frequency/voltage levels. Both algorithms show that migration allows achieving important energy benefits. These benefits are, on average, as much as 17% and 24% for the SOM and MOM algorithms, respectively, over the WF algorithm. An interesting observation, is that global DVS regulators minimize differences among the scheduling strategies for a high number of cores and frequency/voltage levels; showing that, in such a case, SOM achieves many times energy savings close to an idealized scheduler.

The remaining of this paper is structured as follows. Section 2 discusses the related research on energy management and task migration for embedded systems. Section 3 describes the modeled system, including the partitioner and the power-aware scheduler. Section 4 presents the proposed workload partitioning algorithms. Section 5 analyzes energy experimental results. Finally, Section 6 presents some concluding remarks.

2. RELATED WORK

Scheduling in multiprocessor systems can be performed in two main ways depending on the task queue management: *global scheduling*, where a single task queue is shared by all the processors, or *partitioned scheduling*, which uses a private queue for each processor. The former allows by design task migration since all the processors share the same task queue. In the latter case, the scheduling in each processor can be performed by applying well-established uniprocessor theory algorithms such as EDF (Earliest Deadline First) or RMS (Rate Monotonic Scheduling). An example of a modern global scheduling proposal can be found in [13].

In the partitioned scheduling case, research can focus either on the partitioner or the scheduler. Acting in the partitioner, recent works have addressed the energy-aware task allocation problem [14, 15, 9]. For instance, Wei et al. [14] reduce energy consumption by exploiting parallelism of multimedia tasks on a multicore platform combining DVS with switching-off cores. Aydin et al. [15] proposed a algorithm that reserves a subset of processors for the execution of tasks with utilization not exceeding a threshold. Schranzhofer et al. [16] presented a method for allocating tasks to a multiprocessor platform, aimed at minimizing the average power consumption, however, the application is modeled without considering timing constraints. Unlike our work, none of these techniques analyzes the power benefits of task migration among cores.

Some proposals dealing with task migration can be found in the bibliography. Brandenburg et al. [17] evaluate global and partitioned scheduling algorithms in terms of scalability, although power consumption was not investigated. In [18], Zheng divides tasks into fixed and migration tasks, allocating each of the latter category to two cores, so they can migrate from one to another. Unlike this, that paper does not consider dynamic workload changes, instead, all tasks are assumed to arrive at the same instant, so migrations can be scheduled off-line. Seo et al. [7] present a dynamic repartitioning algorithm with migration to balance the workload and reduce consumption. In [19] Brião et al. analyze how migrating soft tasks affects NoC-based MPSoCs in terms of deadline misses and energy consumption. However, the two latter works focus on non-threaded architectures.

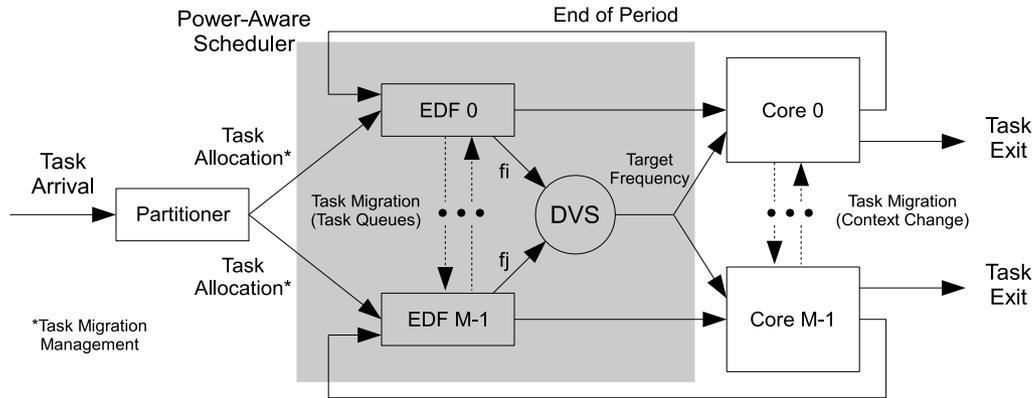


Figure 1. Modeled system.

Regarding the scheduler, in [20] the authors virtualize a simultaneous multithreaded (SMT) processor into multiple single-threaded superscalar processors with the aim of combining high performance with real-time formalism. In order to improve real-time tasks predictability, Cazorla et al. [21] devise an interaction technique between the Operating System (OS) and an SMT processor. Fisher and Baruah [22] derived near-optimal sufficient tests for determining whether a given collection of jobs with precedence constraints can feasibly meet all deadlines upon a specified multiprocessor platform under global EDF scheduling. In [23], the authors propose a methodology for abstracting the total computing power available on a multicore platform by a set of virtual processors, to allocate real-time tasks. The sets of tasks are automatically partitioned into a set of subgraphs that are selected to minimize either the overall bandwidth consumption or the required number of cores. Notice that these works do not tackle energy consumption.

3. SYSTEM MODEL

Figure 1 shows a block diagram of the modeled system. When a task arrives to the system, a partitioner module allocates it into a task queue associated to a given core, which contains the tasks that are ready for execution in that core. These queues are components of the power-aware scheduler that controls a global DVS regulator. In this scheme, the scheduler is in charge of adjusting the working frequency of the cores in order to satisfy the workload requirements.

Processor cores, modelled as an ARM11 MPCore, implement the coarse-grain multithreading paradigm that switches the running thread when a long latency event occurs (i.e., a main memory access). Thus, the running thread issues instructions to execute while another thread is performing the memory access, so overlapping their execution. In the modeled system, the issue slots are always assigned to the thread executing the task with the highest real-time priority. If this thread stalls due to a long latency memory event, then the issue slots are temporarily reassigned until the event is resolved.

3.1. Real-Time Task Behavior

The system workload executes periodic hard real-time tasks. There is no task dependency and each task has its own period of computation. A task can be launched to execute at the beginning of each active period, and it must end its execution before reaching its deadline. The end of the period and the deadline of a task are assumed to be the same for a more tractable scheduling process. There are also some periods where tasks do not execute since they are not active (i.e., inactive periods). In short, a task arrives to the system, executes several times repeatedly, leaves the system, remains out of the system for some periods, and then it enters the system again. This sequence of consecutive active and inactive periods allows modelling real systems with mode changes.

Besides its period and deadline, a task is also characterized by its Worst Case Execution Time (WCET). The task utilization is obtained as $U = \frac{WCET}{Period}$ and is used by several schedulers and partitioners to check whether schedulability of the task set is feasible or not.

3.2. Power-Aware Scheduler

Once a task is allocated to a core, it is inserted into the task queue of that core, where incoming tasks are ordered according to the EDF policy, which prioritizes the tasks with the closest deadlines. Thus, the tasks with the closest deadlines will be the ones mapped into the hardware threads implemented in each core.

The scheduler is also in charge of calculating the required target speed of each core according to the tasks' requirements. In this sense, the EDF scheduler of each core chooses the minimum frequency that fulfills the temporal constraints of its task set in order to minimize power consumption. This information is sent to the global DVS regulator that selects the maximum frequency/voltage level among the requested by the EDF schedulers.

The target frequencies are recalculated only when the workload changes, that is, when a task arrives to and/or leaves the system. In the former case, a higher speed can be required because the workload increases. In the latter, it could happen that a lower frequency could satisfy the deadline requirements of the remaining tasks.

Different frequency values are considered for the power-aware scheduler, based on the frequency levels of a Pentium M [24] which are shown in Table I. This work evaluates the benefits of a DVS with 8, 4 and 2 frequency/voltage levels. The 8L configuration allows the system to work at all the frequencies indicated in the table, whereas the 4L mode permits running tasks at 1700, 1400, 1100 and 600 MHz. The last DVS configuration, referred to as 2L, only supports the extreme frequencies (i.e., 600 and 1700 MHz). In addition, the overhead of changing the frequency/voltage level has been modeled according to a voltage transition rate of $1mv/1\mu s$ [2].

Table I. Frequency (F) vs Power (P).

F[MHz]	1700	1500	1400	1300	1200	1100	900	600
P[Watts]	24.5	24.5	22	22	12	12	7	6

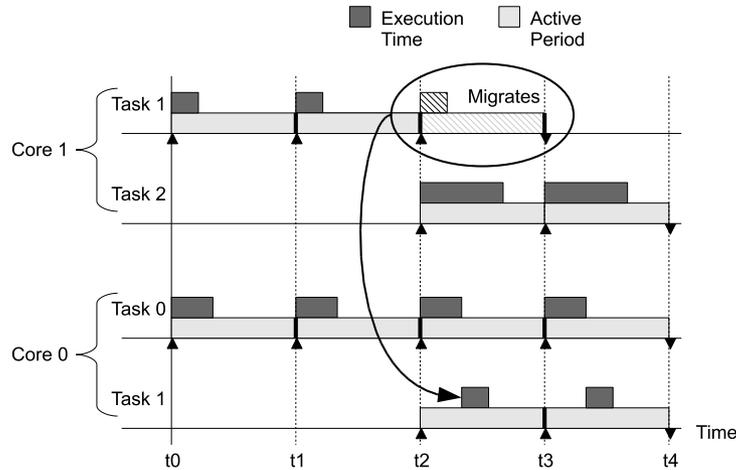


Figure 2. Example of task migrations to balance the system.

4. PROPOSED PARTITIONING HEURISTICS FOR A MULTICORE SYSTEM

There are several partitioning heuristics that can be used to distribute tasks among cores as they arrive to the system. The Worst Fit (WF) partitioning heuristic is considered as one of the best choices in order to balance the workload [9], so yielding to improved energy savings. WF balances the workload by assigning each incoming task to the least loaded core. If more than one task arrives to the system at the same time, WF arranges the incoming tasks in a decreasing utilization order and assigns them to the cores starting with the task with the highest utilization. This algorithm was originally used in partitioned scheduling, and it does not support any task migration among cores by design. In other words, once WF assigns an incoming task to a given core, the task remains in that core until it leaves the system (i.e., it has executed all its active periods). To allow migration, SOM policies are devised in the next subsection.

4.1. Single Option Migration Policies

Figure 2 shows an example of how task migration could improve workload balance. At the beginning of the execution (time t_0), *task 0* and *task 1* are the only tasks assigned to *core 0* and *core 1*, respectively. *Task 0* presents an utilization around 33%, while the utilization of *task 1* is around 25% (i.e., its WCET occupies a quarter of its period). At point t_2 , *task 2*, whose utilization is around 66%, arrives to the system. The WF algorithm would assign it to *core 1* (since it is the least loaded core); leading the system to a high workload imbalance since the global utilization of *core 0* and *core 1* would be 33% and 91%, respectively. This imbalance problem could be solved by allowing task migration. For instance, allowing *task 1* to migrate to *core 0*, would provide a much fair balance (58% in *core 0* versus 66% in *core 1*).

This paper assumes that the running workload dynamically changes at run-time. In this context, the system can mainly become strongly unbalanced when the workload changes, that is, when a task enters or leaves the system, as seen in the previous example. Thus, in the evaluated system migration policies should apply in these points in order to maximize benefits due to migration. For

this purpose, we have devised three policies based on the WF to explore energy benefits: (SOM_{in}), SOM_{out} , and SOM_{in-out} . The first one, SOM_{in} , allows migration only when a new task arrives to the system, SOM_{out} when a task leaves the system, and the last one, SOM_{in-out} , allows migration in both cases. To avoid performing an excessive number of migrations, which could lead to an unacceptable overhead, the number of migrations is limited to only one that can be performed when a task arrives to or leaves the system.

Figure 3 illustrates the devised Migration Attempt (MA) algorithm. This algorithm calculates the imbalance by subtracting the utilization of the least loaded core from the utilization of the most loaded one. This result is divided by two to obtain a theoretical utilization value that represents the amount of work that should migrate to achieve a perfect balance between the two cited cores, and hence, a better global balance. Then, it searches the task in the most loaded core whose utilization is the closest one to this value. Notice that if the utilization of the selected task is not close enough, the migration could yield to a worse imbalance; therefore, the algorithm performs the migration only if it effectively reduces the imbalance.

4.2. Multiple Option Migration Dynamic Partitioner

This subsection presents the Multiple Option Migration (MOM) dynamic partitioner algorithm, which applies both at tasks' arrivals and exits. When a task arrives to the system, MOM selects the target core and performs a migration attempt according to the MA algorithm discussed above. When a task leaves the system MOM checks if a migration attempt would provide energy improvements.

MOM (Figure 4) arranges the tasks arriving to the system in decreasing utilization order. Then it iteratively performs a tentative assignment of the task showing more utilization to each core in order to find which assignment provides the minimum utilization for the most loaded core (U_{min} variable in the figure). Notice that all the possible assignments include a migration attempt according to the MA algorithm discussed above. Finally, the task assignment that provides the best overall balance is applied and the algorithm continues with the next task.

Figure 5 depicts an example where the MOM heuristic improves the behavior of SOM_{in-out} on a task arrival. The SOM_{in-out} allocates the incoming task to core 0 and then performs a migration

```

1:  $imbalance \leftarrow max\_core\_utilization - min\_core\_utilization$ 
2:  $target\_utilization \leftarrow imbalance/2$ 
3:  $minimum\_difference \leftarrow MAX\_VALUE$ 
4: for all  $task$  in  $most\_loaded\_core$  do
5:   if  $|U_{task} - target\_utilization| < minimum\_difference$  then
6:      $minimum\_difference \leftarrow |U_{task} - target\_utilization|$ 
7:      $candidate \leftarrow task$ 
8:   end if
9: end for
10:  $new\_max\_core\_utilization \leftarrow max\_core\_utilization - U_{candidate}$ 
11:  $new\_min\_core\_utilization \leftarrow min\_core\_utilization + U_{candidate}$ 
12:  $new\_imbalance \leftarrow |new\_max\_core\_utilization - new\_min\_core\_utilization|$ 
13: if  $new\_imbalance < imbalance$  then
14:    $migrate(candidate)$ 
15: end if

```

Figure 3. Migration Attempt algorithm.

attempt, but in this case, there is not any possible migration enabling a better workload balance. Thus, the final imbalance becomes 40% (i.e., 90% – 50%). In contrast, when MOM is applied, it also checks the result of allocating the new task to core 1 (arrow labeled as MOM B) and then considering one migration. In this case, the task migration enables a better balance since both cores remain equally loaded with 70% of utilization, which is the distribution selected by MOM.

To sum up, the main difference between SOM_{in-out} and MOM is that the former selects only one core and performs a migration attempt, whereas the proposed heuristic checks different cores, and then chooses the best option in terms of workload balance.

5. EXPERIMENTAL RESULTS

Experimental evaluation has been conducted on Multi2Sim [25], a cycle-by-cycle execution driven simulation framework for evaluating multicore multithreaded processors, which has been extended to model the system described in Section 3. This section evaluates a multicore processor with two, three and four cores, implementing three hardware threads each. Internal core features have been modeled like an ARM11 MPCore based processor, but it has been modified to work as a coarse-grain multithreaded processor with in-order execution, two-instruction issue width, and a 34-cycle memory latency. Regarding the migration overhead, a 10.000 cycles penalty has been assumed [26]. This penalty is applied each time a running context moves its execution to another core.

Since some time is needed to overcome the voltage difference between two different DVS levels, frequency changes are not instantaneous. To model this latency and the power overhead caused by these changes, the worst case for that transition has been assumed. That is, during a frequency

```

1: Algorithm: Multiple Option Migration dynamic partitioner (MOM)
2: Input:  $Task\_set(Task_0, Task_1, \dots, Task_{T-1})$ : task set to be distributed;
3: Input:  $T$ : number of tasks
4: Input:  $Core\_set(Core_0, Core_1, \dots, Core_{M-1})$ : cores in the system
5: Input:  $M$ : number of cores
6: Input/Output:  $Tasks_0, Tasks_1, \dots, Tasks_{M-1}$ : tasks sets assigned to the different M cores
7: while  $Task\_set$  is not empty do
8:    $target\_task \leftarrow Task_i : (Task_i) \geq MAX(U(Task_0), U(Task_1), \dots, U(Task_{T-1}))$ 
9:    $U_{min} \leftarrow \infty$ 
10:   $initial\_task\_assignment = Tasks_0, Tasks_1, \dots, Tasks_{M-1}$ 
11:  for all  $target\_core$  in  $Core\_set$  do
12:     $Tasks_{target\_core} \leftarrow Tasks_{target\_core} \cup \{target\_task\}$ 
13:     $Migration\_Attempt()$ 
14:    if  $U_{min} > MAX(U(Core_0), U(Core_1), \dots, U(Core_{M-1}))$  then
15:       $U_{min} \leftarrow MAX(U(Core_0), U(Core_1), \dots, U(Core_{M-1}))$ 
16:       $best\_task\_assignment \leftarrow Tasks_0, Tasks_1, \dots, Tasks_{M-1}$ 
17:    end if
18:     $Tasks_0, Tasks_1, \dots, Tasks_{M-1} \leftarrow initial\_task\_assignment$ 
19:  end for
20:   $Tasks_0, Tasks_1, \dots, Tasks_{M-1} \leftarrow best\_task\_assignment$ 
21: end while

```

Figure 4. Multiple Option Migration dynamic partitioner algorithm.

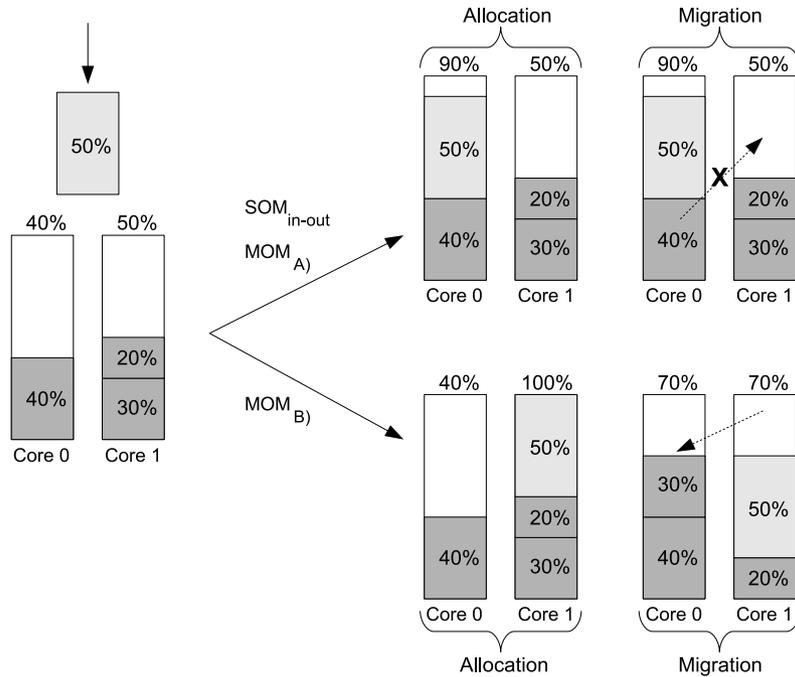


Figure 5. SOM_{in-out} vs MOM working example.

transition the speed of the lowest frequency and the power consumption of the highest one are taken into account.

Table II shows the benchmarks from [27] that have been used to prepare real-time workload mixes (a mix number with an asterisk means that the benchmark is used in the mix more than once). Each mix is composed of a set of benchmarks whose number ranges from 7 to 34, running concurrently depending on the number of cores executed. Mixes 1, 2 and 3 are executed in a 2-core system, mixes 4, 5 and 6 in a system with three cores, and mixes 7, 8 and 9 in a 4-core system. These mixes have been designed taking into account aspects such as task utilization, task periodicity, and the sequence of active and inactive periods. Task periods range from 100.000 to 18.000.000 cycles, the number of times that a task arrives to and leaves the system from 1 to 21, and the consecutive number of active periods of a task from 1 to 70. The global system utilization varies in a single execution from 35% to 95% in order to test the algorithms behavior across a wide range of situations. In addition, all results are presented and analyzed for a system implementing two, four and eight voltage levels.

5.1. Impact of Applying Migrations at Specific Points of Time

This section analyzes the three devised Single Option Migration variants (SOM_{in} , SOM_{out} and SOM_{in-out}). The main goal is to identify the best points of time to carry out migrations. Figure 6 shows the relative energy consumption compared to the energy consumed by the system working always at the maximum speed, for different benchmark mixes, DVS configurations, and number of cores. The results are obtained by multiplying the number of cycles working at each frequency by the energy required per cycle at that frequency.

Table II. Benchmark description and mixes. Legend: * the benchmark appears more than once in the mix.

Name	Function Description	Mix
bs	Binary search for a 15-element array	1, 3, 4*, 9*
bsort100	Bubblesort program	3
cnt	Counts non-negative numbers in a matrix	2, 3, 4, 7*
compress	Data compression program	2, 3, 4, 7*
cover	Program for testing many paths	5*, 8*
crc	Cyclic redundancy check on 40-byte data	7*
duff	Copy 43-byte array	3, 4*
edn	FIR filter calculations	7*
expint	Series expansion for integral function	2, 3, 4*, 7*
fac	Factorial of a number	1, 2, 3, 4*, 7*, 9*
fdct	Fast Discrete Cosine Transform	5, 8*
fft1	1024-point Fast Fourier Transform	2, 3, 4*, 7*
fibcall	Simple iterative Fibonacci calculation	1, 3, 4*, 6*, 9*
fir	Finite impulse response filter	5, 8*
insertsort	Insertion sort on a reversed array of size 10	3, 4*
janne_complex	Nested loop program	1, 2, 5*, 7*, 8*, 9*
jfdctint	Discrete-cosine transformation	2, 5*, 7*, 8*
lcdnum	Read ten values, output half to LCD	1, 3, 4*, 6*, 9*
loop3	Function with diverse loops	3, 6*
ludcmp	LU decomposition algorithm	2, 5*, 7*, 8*
minmax	Minimum and maximum functions	5*, 6*, 8*
minver	Inversion of floating point matrix	3, 4*
ns	Search in a multi-dimensional array	3
nsichneu	Simulate an extended Petri Net	5*, 8*
qsort-exam	Non-recursive version of quick sort algorithm	5*, 8*
qurt	Root computation of quadratic equations	2, 5*, 7*, 8*
select	Nth largest number in a floating point array	5*, 6*, 8*
sqrt	Square root function	1, 5*, 6*, 8*, 9*
statemate	Automatically generated code	1, 3, 4*, 9*

As it can be observed in the results of the 2-core system (Figure 6(a)), migration can provide important energy savings with respect to no migration (WF). For instance, for mix 2 in the 4L case with task migration, both when a task arrives to and leaves the system, the energy consumption can be reduced by up to 23.27% compared to the execution without migration.

An interesting observation is that, in some mixes, the SOM_{in} variant consumes more power than the classical WF algorithm with no migration. For example, in the 3-core system (Figure 6(b)) allowing migrations only at tasks' arrivals turns out in harmful effects for mix 4 in terms of power consumption, where SOM_{in} consumes 12.27% more energy than WF for 4L configuration. The reason is related to the fraction of time length that the system is *controlled* by the partitioning algorithm. That is, the SOM_{in} partitioning heuristic only applies at tasks' arrivals. Therefore, as soon as a task leaves the system, the workload imbalance will rise since SOM_{in} does not apply on such events.

Figure 7 illustrates an example. At *time* t_0 tasks T_1 , T_2 , and T_3 arrive to the system, and the scheduler selects the frequency/voltage level that best fits the workload requirements. Lets assume that the workload is perfectly balanced in a 2-core system. Then at *time* t_1 , task T_1 leaves the

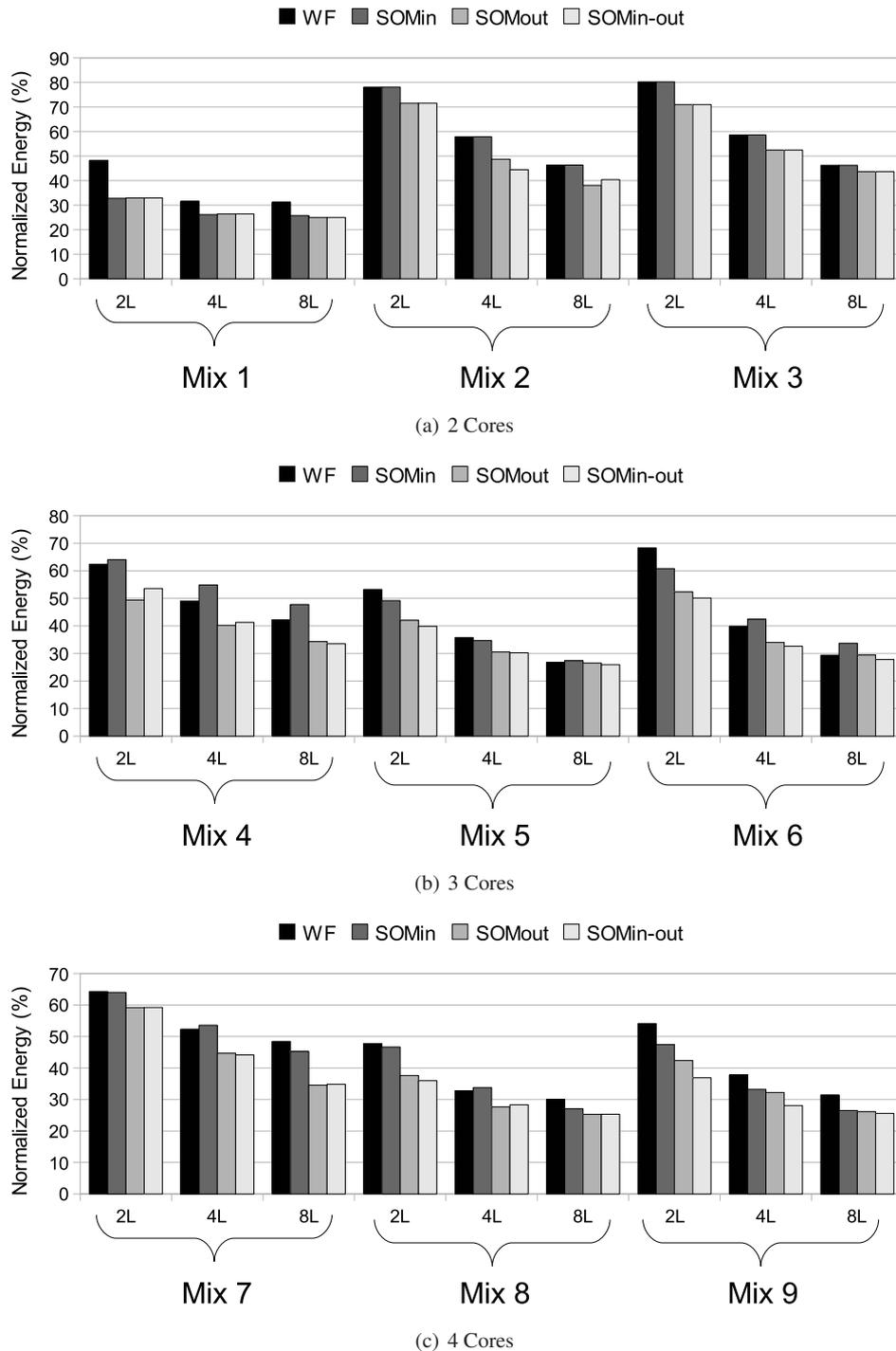


Figure 6. Single Option Migration variants comparison for different DVS levels and number of cores.

system, so workload imbalance will rise (dashed area), in algorithms such WF or SOM_{in} where migration is not performed, so yielding to energy wasting. Notice that this area is *uncontrolled* since the set of tasks running has changed. On the contrary, the *controlled* time periods are those where

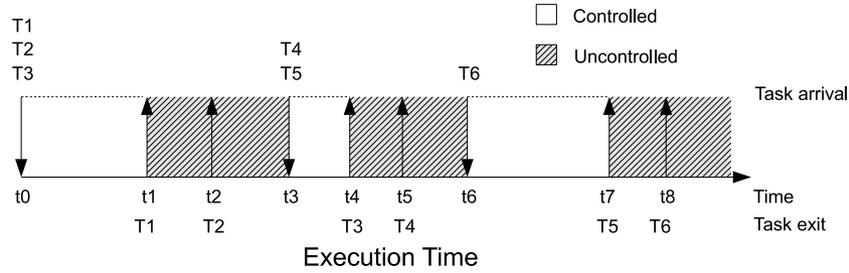


Figure 7. Effective action of the SOM_{in} partitioning algorithm.

the set of tasks running matches the set used to perform the scheduling actions. Moreover, further imbalance would rise if the next task $T2$ leaves the system from the same core. This imbalance will remain until the algorithm applies, which happen only on tasks' arrivals in WF and SOM_{in} (in $t3$). This drawback is solved in the algorithms which allow task migration at such points like SOM_{out} , SOM_{in-out} and MOM. Table III shows which actions are performed by the different algorithms both when a task arrives to and leaves the system.

The longer the algorithm controls the running workload, the better the workload balance. Consequently, the frequency levels requested by the different cores will be similar, so avoiding energy wasting. Figure 8 shows, for mix 4, in a 3-core system with 8 DVS levels, the difference among frequencies required by the cores along the execution time (in percentage). For instance label 0 means that both cores require the same frequency and label 2 means that the core with less frequency/voltage requirements requested level i to the DVS regulator, while the core with the maximum requirements requested level $i+2$. This figure explains the curious behavior identified above, where SOM_{in} performed worse than WF. As observed, both partitioners yield the system to spend a similar amount of time with all the cores requiring a similar speed (i.e., with a difference less or equal than 1 level). Nevertheless, the main reason why SOM_{in} consumes more power than WF is that, in this mix, there is a significant amount of time where the difference in speed required by the cores in SOM_{in} is 3 and 4 levels, while in WF most of this time the difference is only 2 levels. Notice that SOM_{out} and SOM_{in-out} balance the workload in a better way (area associated to label 0 is much longer) than WF and SOM_{in} , the reason is due to the former control the system both at tasks' arrivals and exits.

Another interesting remark is that if the system implements more DVS frequency levels, then more energy savings can be potentially obtained since the system can select a frequency closer to

Table III. Algorithms action on workload changes.

Algorithm	Task Arrival	Task Exit
WF	WF	-
SOM_{in}	WF, Task Migration	-
SOM_{out}	WF	Task Migration
SOM_{in-out}	WF, Task Migration	Task Migration
MOM	MOM	Task Migration

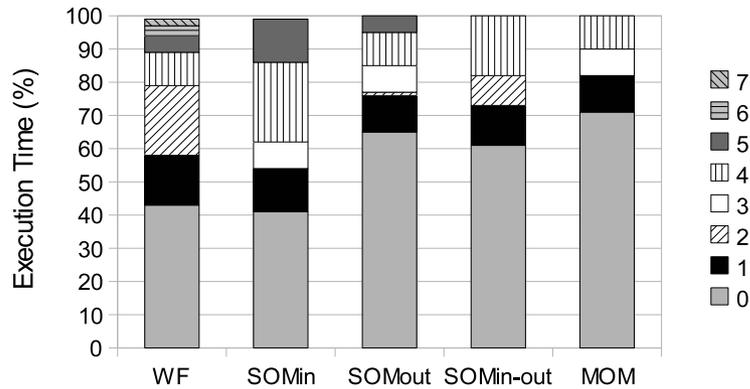


Figure 8. Differences of the required frequencies.

the optimal estimated by the scheduler. However, despite this fact, in some cases energy benefits due to migration in a system with few frequency levels can reach or even surpass the benefits of having more levels without migration. For example, the energy consumption of SOM_{in-out} for mix 2 in the 4L 2-core system is around 44% the consumption of the baseline, whereas the same value of WF in the 8L system is 46%.

Finally, it can be noticed that the system behaves in a similar way regardless of the number of cores, that is, the benefits of migration that are observed in systems with two or three cores are also similar in a system with four cores, as shown in Figure 6(c). This fact makes the proposal a good candidate for commercial systems attending to the current industry trend of increasing the number of cores.

5.2. Comparing MOM Versus SOM Variants

This section analyzes the energy improvements of the proposed MOM algorithm over the SOM variants algorithms. For comparison purposes the best SOM variant (SOM_{in-out}), on average, and a theoretical threshold have been also included in the plots. This theoretical threshold is a value that represents the maximum energy savings that can be achieved in a system where the number of task migrations is not limited, they have no cost, and they can be performed at any point of the execution time. That is, a system with perfect task balancing and without penalties due to migration. Figure 9 shows the energy results for two, three and four cores, normalized with respect to the energy consumed by the system working always at the maximum speed.

Results show that, regardless of the number of cores, the mix, and the number of frequency levels, MOM saves more energy than SOM_{in-out} . For example, when running mix 3 in the 2L 2-core system, MOM consumes 60.17% and 68.01% of the energy consumed by WF and by SOM_{in-out} , respectively. The reason is that MOM enables the cores of the system to work at a similar frequency for longer than any SOM variant. This can be also observed in the example of Figure 8. Comparing the working behavior of MOM with SOM_{in-out} it can be appreciated that both algorithms perform the same action when a task leaves the system (see Table III). Therefore, differences in benefits between them come from applying the algorithm at tasks' arrivals. The reason is that SOM_{in-out}

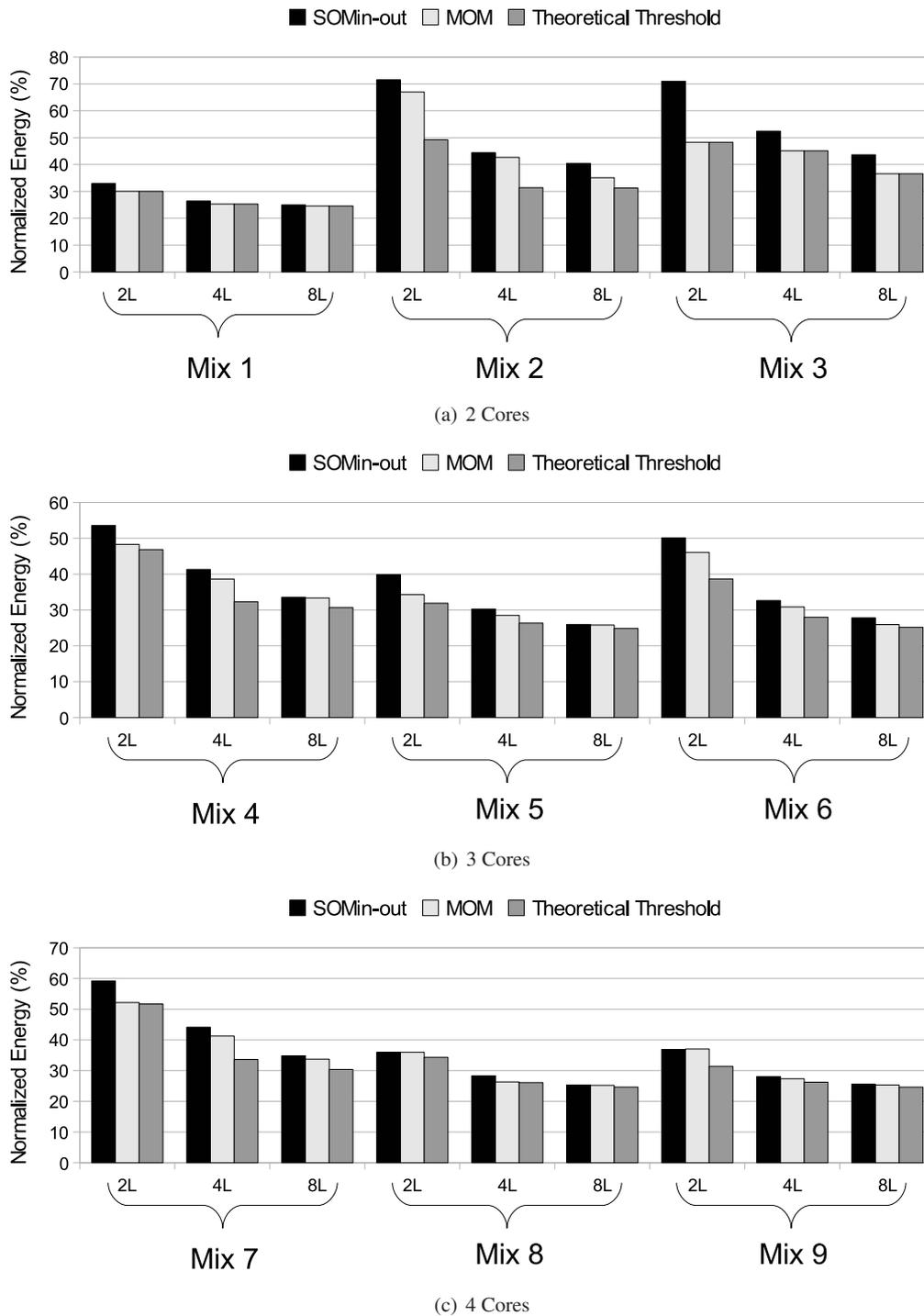


Figure 9. SOM_{in-out} versus MOM for different DVS levels and number of cores.

first allocates the incoming task and then makes one migration attempt, whereas MOM checks for each core which combination of task-to-core allocation plus a single migration attempt would achieve a better workload balance. Thus, MOM examines a wider range of possible distributions.

Moreover, in some mixes (e.g., mix 3) MOM results are very close to the energy savings of the theoretical threshold. However, if the utilizations of the tasks in a given mix widely differ among them, and depending on how run-time conditions evolve, the results of any practical scheduler may be far from the theoretical threshold (e.g., mix 7, mix 9). The standard deviation of the task utilization (see Table IV) in these mixes is relatively high since a few tasks have a huge utilization. This fact prevents SOM and MOM from achieving a perfect balancing in some scenarios, as done by the theoretical threshold. Notice that mix 1 for a 2-core system also presents a high standard deviation value, but in this case it is due to a single task with much higher utilization. On the other hand, in mix 3 most tasks present similar utilizations within a limited range (10%-17%). Thus, it is more feasible that practical schedulers can obtain a perfect balancing.

Finally, as the number of cores and voltage levels increases (4 cores), a Single Option Migration algorithm is enough to achieve important energy savings, although MOM can slightly improve those results. Moreover, these results fall close to the theoretical maximum. Thus, in this scenario, a possible choice to enhance energy savings is to change the voltage regulator domain (i.e., to implement several regulators, each one shared by a subset of the cores).

6. CONCLUSIONS

Workload balancing has been proved to be an efficient power technique in multicore systems. Unfortunately, unexpected workload imbalances can rise at run-time provided that the workload changes dynamically since new tasks arrive to or leave the system. To palliate this shortcoming, this paper has analyzed the impact on energy consumption due to scheduling strategies in a multicore embedded system implementing DVS.

Two power-aware schedulers working with real-time constraints, namely SOM and MOM have been devised, which check only one target core or the optimal core before performing a migration, respectively. To prevent excessive overhead, task migration has been strategically applied at three specific execution points of time where the workload changes: at tasks' arrivals, at tasks' exits, and in both cases. Three variants of SOM algorithm are devised depending on the point of time the algorithm applies.

Experimental evaluation has been performed using sets of mixes of real-time benchmarks executed on a modeled ARM11 MPCore processor. A first observation is that applying the algorithm at tasks' exits achieves better energy savings than applying it only at tasks' arrivals, but the highest benefits are obtained when the algorithm is applied in both cases. On the other hand, MOM performs in general better than SOM, however as the number of cores and frequency/voltage levels increases, the differences among energy benefits are reduced. Results show that task migration allows the proposed schedulers to achieve important energy benefits over the WF. These benefits are, on

Table IV. Average and Standard Deviation of Task Utilization.

Mix	1	2	3	4	5	6	7	8	9
Average	30.54	24.29	15.32	14.96	14.86	20.28	19.94	13.33	16.06
Standar Deviation	13.86	8.12	5.29	3.39	2.76	4.44	12.07	3.64	10.46

average, by 17% and 24% over the WF, for the SOM and MOM, respectively. Moreover, in some cases MOM's benefits are up to 40%.

This paper has shown how task migration combined with DVS can allow important energy savings. Thus, benefits come from both techniques. Analyzing the results one can notice that migration is a powerful technique since it allows reducing energy consumption compared to a system with more voltage levels without migration.

A final remark is that improving the workload balance by supporting task migration, not only energy savings can be enhanced, but since the utilization of the most loaded core is also reduced, then also a wider set of tasks could be scheduled.

References

1. Donald J, Martonosi M. Techniques for Multicore Thermal Management: Classification and New Exploration. *Proceedings of the 33rd Annual International Symposium on Computer Architecture*, IEEE Computer Society: Boston, MA, USA, 2006; 78–88.
2. Wu Q, Martonosi M, Clark DW, Reddi VJ, Connors D, Wu Y, Lee J, Brooks D. A Dynamic Compilation Framework for Controlling Microprocessor Energy and Performance. *Proceedings of the 38th Annual IEEE/ACM International Symposium on Microarchitecture*, IEEE Computer Society: Barcelona, Spain, 2005; 271–282.
3. McNairy C, Bhatia R. Montecito: A Dual-Core, Dual-Thread Itanium Processor. *IEEE Micro* 2005; **25**(2):10–20.
4. Kalla R, Sinharoy B, Tendler J. IBM Power5 Chip: A Dual-Core Multithreaded Processor. *IEEE Micro* 2004; **24**(2):40–47.
5. Shah A. *Arm plans to add multithreading to chip design*. ITworld 2010. [Online]. Available: <http://www.itworld.com/hardware/122383/arm-plans-add-multithreading-chip-design>.
6. Hung C, Chen J, Kuo T. Energy-Efficient Real-Time Task Scheduling for a DVS System with a Non-DVS Processing Element. *Proceedings of the 27th Real-Time Systems Symposium*, IEEE Computer Society: Rio de Janeiro, Brazil, 2006; 303–312.
7. Seo E, Jeong J, Park S, Lee J. Energy Efficient Scheduling of Real-Time Tasks on Multicore Processors. *IEEE Transactions on Parallel and Distributed Systems* 2008; **19**(11):1540–1552.
8. March J, Sahuquillo J, Hassan H, Petit S, Duato J. A New Energy-Aware Dynamic Task Set Partitioning Algorithm for Soft and Hard Embedded Real-Time Systems. *The Computer Journal* 2011; [Online] 8 February 2011, ISSN: 1460–2067.
9. AlEnawy TA, Aydin H. Energy-Aware Task Allocation for Rate Monotonic Scheduling. *Proceedings of the 11th Real Time on Embedded Technology and Applications Symposium*, IEEE Computer Society: San Francisco, CA, USA, 2005; 213–223.
10. INTEL Corp., Santa Clara, CA, USA. *Intel Atom Processor Microarchitecture*. [Online]. Available: www.intel.com/.
11. Marvell Semiconductor, Inc., Santa Clara, CA, USA. *Marvell ARMADATM 628*. [Online]. Available: http://www.marvell.com/company/press_kit/assets/Marvell_ARMADA_628_Release_FINAL3.pdf.
12. Hirata K, Goodacre J. ARM MPCore; The streamlined and scalable ARM11 processor core. *Proceedings of the Conference on Asia South Pacific Design Automation*, IEEE Computer Society: Yokohama, Japan, 2007; 747–748.
13. Kato S, Yamasaki N. Global EDF-based Scheduling with Efficient Priority Promotion. *Proceedings of the 14th International Conference on Embedded and Real-Time Computing Systems and Applications*, IEEE Computer Society: Kaohsiung, Taiwan, 2008; 197–206.
14. Wei Y, Yang C, Kuo T, Hung S. Energy-Efficient Real-Time Scheduling of Multimedia Tasks on Multi-Core Processors. *Proceedings of the 25th Symposium on Applied Computing*, ACM: Sierre, Switzerland, 2010; 258–262.
15. Aydin H, Yang Q. Energy-Aware Partitioning for Multiprocessor Real-Time Systems. *Proceedings of the 17th International Parallel and Distributed Processing Symposium, Workshop on Parallel and Distributed Real-Time Systems*, IEEE Computer Society: Nice, France, 2003; 113.
16. Schranzhofer A, Chen JJ, Thiele L. Dynamic power-aware mapping of applications onto heterogeneous MPSoC platforms. *IEEE Transactions on Industrial Informatics* 2010; **6**(4):692–707.
17. Brandenburg BB, Calandrino JM, Anderson JH. On the Scalability of Real-Time Scheduling Algorithms on Multicore Platforms: A Case Study. *Proceedings of the 29th Real-Time Systems Symposium*, IEEE Computer Society: Barcelona, Spain, 2008; 157–169.

18. Zheng L. A Task Migration Constrained Energy-Efficient Scheduling Algorithm for Multiprocessor Real-time Systems. *Proceedings of the International Conference on Wireless Communications, Networking and Mobile Computing*, IEEE Computer Society: Shanghai, China, 2007; 3055–3058.
19. Brião E, Barcelos D, Wronski F, Wagner FR. Impact of Task Migration in NoC-based MPSoCs for Soft Real-time Applications. *Proceedings of the International Conference on VLSI*, IEEE Computer Society: Atlanta, GA, USA, 2007; 296–299.
20. El-Haj-Mahmoud A, AAL-Zawawi, Anantaraman A, Rotenberg E. Virtual Multiprocessor: An Analyzable, High-Performance Architecture for Real-Time Computing. *Proceedings of the International Conference on Compilers, Architectures and Synthesis for Embedded Systems*, ACM Press: San Francisco, CA, USA, 2005; 213–224.
21. Cazorla F, Knijnenburg P, Sakellariou R, Fernández E, Ramirez A, Valero M. Predictable Performance in SMT Processors: Synergy between the OS and SMTs. *IEEE Transactions on Computers* 2006; **55**(7):785–799.
22. Fisher N, Baruah S. The feasibility of general task systems with precedence constraints on multiprocessor platforms. *Real-Time Systems* 2009; **41**(1):1–26.
23. Buttazzo G, Bini E, Wu Y. Partitioning Real-Time Applications Over Multicore Reservations. *IEEE Transactions on Industrial Informatics* 2011; **7**(2):302–315.
24. INTEL Corp., Santa Clara, CA, USA. *Intel Pentium M Processor Datasheet* 2004. [Online]. Available: <http://download.intel.com/support/processors/mobile/pm/sb/25261203.pdf>.
25. Ubal R, Sahuquillo J, Petit S, López P. Multi2Sim: A Simulation Framework to Evaluate Multicore-Multithreaded Processors. *Proceedings of the 19th International Symposium on Computer Architecture and High Performance Computing*, IEEE Computer Society: Gramado, RS, Brazil, 2007; 62–68.
26. Chaparro P, González J, Magklis G, Cai Q, González A. Understanding the Thermal Implications of Multi-Core Architectures. *IEEE Transactions on Parallel and Distributed Systems* 2007; **18**(8):1055–1065.
27. Real-Time Bench, Vasteras, Sweden. *WCET Analysis Project. WCET Benchmark Programs* 2006. [Online]. Available: <http://www.mrtc.mdh.se/projects/wcet/>.