# Energy-Efficient Thermal-Aware Scheduling for RT Tasks Using TCPN

## L. Rubio-Anguiano \* G. Desirena-López \* A. Ramírez-Treviño \* J.L. Briz \*\*

\* CINVESTAV-IPN Unidad Guadalajara, Av. del Bosque 1145, CP 45019, Zapopan, Jalisco, Mexico \*\* DIIS/I3A Univ. de Zaragoza, María de Luna 1 - 50018 Zaragoza, España.

#### Abstract:

This work leverages TCPNs to design an energy-efficient, thermal-aware real-time scheduler for a multiprocessor system that normally runs in a low state energy at maximum system utilization but its capable of increasing the clock frequency to serve aperiodic tasks, optimizing energy, and honoring temporal and thermal constraints. An off-line stage computes the minimum frequency required to run the periodic tasks at maximum CPU utilization, the proportion of each task's job to be run on each CPU, the maximum clock frequency that keeps temperature under a limit, and the available cycles (slack) with respect to the system with minimum frequency. Then, a Zero-Laxity online scheduler dispatches the periodic tasks according to the offline calculation. Upon the arrival of aperiodic tasks, it increases clock frequency in such a way that all periodic and aperiodic tasks are properly executed. Thermal and temporal requirements are always guaranteed, and energy consumption is minimized.

Keywords: Real-Time scheduling, Timed Continuous Petri Nets, Modelling.

## 1. INTRODUCTION

There is a great interest in using multicore processors in embedded real-time systems. Multicores reduce the Size, Weight and Power (SWaP) requirements in avionics, have evident advantages in automotive electronics and satellite systems, and are already common in consumer devices. Limited power, battery lifespan and thermal restrictions require careful energy management, for which real-time schedulers can leverage power management mechanisms provided by current MPSoCs, such as dynamic voltage and frequency scaling (DVFS). Many DVFS schedulers, thermal-aware schedulers and schedulers that minimize energy consumption have been proposed over the past few years (Kong et al. (2014), Schor et al. (2012), Ahmed et al. (2016), Hettiarachchi et al. (2014)).

In this work we present a minimal energy, thermal-aware real-time scheduler for a multiprocessor system. We consider a set of periodic independent tasks subject to temporal, thermal and energy restrictions, modeled by means of a Timed Continuos Petri Net (TCPN). Aperiodic tasks arrive asynchronously, and should only be served when every periodic task is guaranteed to meet its deadline. An off-line stage first computes the minimum frequency required to run the periodic task set at maximum CPU utilization. By means of a Linear Programming Problem (LPP) it finds a maximum frequency  $(F^+)$  at which the system can run subject to temporal and thermal constraints. The final offline step leverages a deadline partitioning approach (Funk et al. (2011)) to compute a schedule by calculating the proportion of each task that must be executed on each CPU during the interval. In this way, the execution of this

schedule consumes minimum energy and meets temporal and thermal constrain. In a second on-line stage, a Zero Laxity scheduler (Davis and Burns (2011b)) performs the schedule computed off-line for the periodic task set. Upon arrival, aperiodic tasks are serviced if a controller can increase the CPU clock frequency to achieve the proper slack, or are rejected otherwise. The frequency is always kept as low as possible to meet the temporal constraints and minimize energy consumption. As far as we know, this is the first solution based on TCPNs to accommodate aperiodic real-time tasks while preserving temporal and thermal constraints, minimizing energy. Our preliminary results, leveraging a deadline partitioning scheme, yield a good compromise in terms of optimal utilization of CPU, context switches and migrations.

Section 2 provides basic definitions and concepts on Timed Continuous Petri Nets (TCPN), the formal model used in this work to represent tasks, CPUs, energy consumption and temporal and thermal behavior. Section 3 formulates the scheduling problem. Section 4 explains the system model. The off-line schedule calculation is explained in Section 5, whereas Section 6 presents the the on-line scheduler. Section 7 shows the procedure to deal with aperiodic aperiodic tasks. Section 8 shows some examples, and Section 9 summarizes conclusions and further work.

## 2. PERI NETS BACKGROUND

This section provides basic definitions and concepts on Timed Continuous Petri Nets (TCPN), the formal model used in this work to represent tasks, CPUs, energy consumption, temporal and thermal behavior. For a deeper

insight on Petri Nets see Silva and Recalde (2007), David and Alla (2008), Silva et al. (2011).

Definition 2.1. A Petri net (PN) is a 4-tuple  $N = (P, T, \mathbf{Pre}, \mathbf{Post})$  where  $P = \{p_1, ..., p_{|P|}\}$  and  $T = \{t_1, ..., t_{|T|}\}$  are finite disjoint sets of places and transitions. **Pre** and **Post** are  $|P| \times |T|$  **Pre**– and **Post**– incidence matrices, where  $\mathbf{Pre}[i, j] > 0$  (resp.  $\mathbf{Post}[i, j] > 0$ ) if there is an arc going from  $p_i$  to  $t_j$  (or going from  $t_j$  to  $p_i$ ),  $\mathbf{Pre}[i, j] = 0$  (or  $\mathbf{Post}[i, j] = 0$ ) otherwise.

Definition 2.2. A continuous Petri net (ContPN) is a pair ContPN =  $(N, \mathbf{m}_0)$  where  $N = (P, T, \mathbf{Pre}, \mathbf{Post})$  is a PN (PN) and  $\mathbf{m}_0 \in \{\mathbb{R}^+ \cup 0\}^{|P|}$  is the initial marking.

A transition  $t_i$  is enabled at  $\boldsymbol{m}$  if  $\forall p_j \in \boldsymbol{t}_i, \boldsymbol{m}[p_j] > 0$ ; and its enabling degree is defined as  $enab(t_i, \boldsymbol{m}) = \min_{p_j \in \boldsymbol{t}_i} \frac{\boldsymbol{m}[p_j]}{\boldsymbol{Pre}[p_j, t_i]}$ . Firing  $t_i$  in a certain amount  $\alpha \leq enab(t_i, \boldsymbol{m})$  yields a marking  $\boldsymbol{m'} = \boldsymbol{m} + \alpha \boldsymbol{C}[P, t_i]$ , where  $\boldsymbol{C} = \boldsymbol{Post} - \boldsymbol{Pre}$ . If  $\boldsymbol{m}$  is reachable from  $\boldsymbol{m}_0$  by firing the finite sequence  $\boldsymbol{\sigma}$  of enabled transitions, then  $\boldsymbol{m} = \boldsymbol{m}_0 + \boldsymbol{C} \boldsymbol{\vec{\sigma}}$  is named the fundamental equation where  $\boldsymbol{\vec{\sigma}} \in \{\mathbb{R}^+ \cup 0\}^{|T|}$  is the firing count vector, i.e  $\boldsymbol{\vec{\sigma}}[j]$  is the cumulative amount of firings of  $t_j$  in the sequence  $\boldsymbol{\sigma}$ .

Definition 2.3. A timed continuous PN (*TCPN*) is a time-driven continuous-state system described by the tuple  $(N, \lambda, m_0)$  where  $(N, m_0)$  is a continuous PN and the vector  $\lambda \in \{\mathbb{R}^+ \cup 0\}^{|T|}$  represents the transitions rates determining the temporal evolution of the system. Transitions fire according to a certain speed, which is generally a function of the transition rates and the current marking. Such function depends on the semantics associated to the transitions. Under infinite server semantics, the flow through a transition t (or t firing speed, denoted as f(m)) is the product of the rate,  $\lambda_i$ , and  $enab(t_i, m)$ , the instantaneous enabling of the transition, i.e.,  $f_i(m) = \lambda_i enab(t_i, m)$ .

The firing rate matrix is defined by  $\mathbf{\Lambda} = diag(\lambda_1, ..., \lambda_{|T|})$ . For the flow to be well defined, every continuous transition must have at least one input place, so we assume  $\forall t \in$  $T, |\bullet t| \geq 1$ . The "min" in the above definition leads to the concept of configuration. A configuration of a *TCPN* at  $\boldsymbol{m}$  is a set of (p, t) arcs describing the effective flow of each transition, and say that  $p_i$  constrains  $t_j$  for each arc  $(p_i, t_j)$ in the configuration. A configuration matrix is defined for each configuration as follows:

$$\boldsymbol{\Pi}_{\boldsymbol{j},\boldsymbol{i}}(\boldsymbol{m}) = \begin{cases} \frac{1}{\boldsymbol{P}\boldsymbol{r}\boldsymbol{e}[i,j]} & \text{if } p_i \text{ is constraining } t_j \\ 0 & \text{otherwise} \end{cases}$$
(1)

 $f(m) = \Lambda \Pi(m)m$  is the vectorial form of the flow of a transition. The following fundamental equation describes the dynamic behaviour of a *PN* system:

$$\mathbf{m} = Cf(m) = C\Lambda\Pi(m)m$$
 (2)

A control action can be applied to (2) by adding a term  $\boldsymbol{u}$  to every transition  $t_i$  such that  $0 \leq u_i \leq f_i$ , indicating that its flow can be reduced. Thus, the controlled flow of transition  $t_i$  becomes  $w_i = f_i - u_i$  and the forced state equation is:  $\boldsymbol{m} = \boldsymbol{C}[\boldsymbol{f} - \boldsymbol{u}] = \boldsymbol{C}\boldsymbol{w}$ .

## 3. PROBLEM DEFINITION

Definition 3.1. Let  $\mathcal{T} = \{\tau_1, ..., \tau_n\}$  be a set of n independent periodic tasks. Each task is identified by the  $3 - tuple \ \tau_i = (cc_i, d_i, \omega_i)$ , where  $cc_i$  is the worst-case execution in CPU cycles,  $\omega_i$  the period and  $d_i$  is the relative implicit deadline  $(d_i = \omega_i)$  (Baruah et al. (2015)). Let  $\mathcal{P} = \{CPU_1, \ldots, CPU_m\}$  be a set of *m* identical processors with an homogeneous clock frequency  $F \in \mathcal{F} =$  $[F_1, ..., F_{max}]$ . We assume that all task parameters, including task period and CPU cycles are integers and that any task can be preempted at any time. The hyper-period is defined as the period equal to the least common multiple of periods  $H = lcm(\omega_1, \omega_2, \dots, \omega_n)$  of the *n* periodic tasks. A task  $\tau_i$  executed on a processor at its maximum frequency  $F_{max}$ , requires  $c_i = \frac{cc_i}{F_{max}}$  processor time at every  $\omega_i$ interval. The system utilization is defined as the fraction of time during which the processor is busy running the task set i.e.,  $U = \sum_{i=1}^{n} \frac{c_i}{\omega_i}$ . The CPU utilization must be computed and should be less or equal to the number of processors i.e.,  $U \leq m$  (Baruah et al. (1996)).

**Problem** 3.1. Minimum Energy Thermal Aware Real Time Scheduler (*METARTS*). Given the sets of tasks  $\mathcal{T}$  and CPUs  $\mathcal{P}$ , the *METARTS* problem consists in designing an algorithm to allocate within the hyperperiod H the tasks in  $\mathcal{T}$  to the m identical CPUs in such a way that the deadlines for  $\mathcal{T}$  are always satisfied and the CPU temperatures are kept always below a given temperature bound  $T_{max}$  and the consumed energy is minimum.

In addition to previous problem, aperiodic tasks, that must be allocated to CPUs arrive to the system.

Definition 3.2. Aperiodic Task. Let  $\mathcal{T}^a = \{\tau_1^a, \ldots, \tau_z^a\}$  be an independent aperiodic task set. Each  $\tau_i^a$  is defined as a 3-tuple  $(cc_i^a, d_i^a, r_i^a)$  where  $cc_i^a$  (required CPU cycles),  $d_i^a$  (its deadline) are known a priory, and  $r_i^a$  is the arrival time, which is unknown.

#### 4. SYSTEM MODEL

The CPU thermal activity and the tasks fluid allocation to CPUs happens in a continuous space and honor a set of differential equations. The use of a TCPN allows modelling this continuous activity while providing data about the system's state at any moment. This section summarizes the TCPN model of the system (CPUs, tasks, thermal behavior, energy). For details on the thermal model we refer to Desirena-Lopez et al. (2014). The task model is based on Desirena-Lopez et al. (2016). Fig. 1 details a part of the model corresponding to a one task ( $\tau_1$ ) and one CPU ( $CPU_1$ ).

Task and CPU Model.- The places  $p_i^{\omega}$ ,  $p_i^{cc}$  and  $t_i^{\omega}$ , together with their arcs, model the i - th task. The places  $p_{1,j}^{busy}, ..., p_{n,j}^{busy}, p_{1,j}^{exec}, ..., p_{n,j}^{stac}, p_j^{idle}$ , and transitions  $t_{1,j}^{alloc}, ..., t_{n,j}^{alloc}$  and  $t_{1,j}^{exec}$ ,  $..., t_{n,j}^{exec}$  together with their arcs, model the j - th CPU. The weighted arc  $cc_i$  corresponds to the duration of  $\tau_i$  (in CPU cycles). The CPU cycles required to run a task's job are stored in place  $p_i^{cc}$ . The transitions  $t_{i,j}^{exec}$  represent the execution on the  $CPU_j$ . The current marking of places  $p_{i,j}^{exec}$  ( $m_{i,j}^{exec}$ ) represents the CPU worst-case execution time of task  $\tau_i$  currently



Fig. 1. Equivalent TCPN model for a multiprocessor system.

executed at a frequency F. Places  $p_j^{idle}$  and  $p_{i,j}^{busy}$  represent the idle state and the busy state of the processor. The marking of place  $p_j^{idle}$  models the available CPU cycles (throughput capacity). The initial marking at  $p_j^{idle}$  is set to 1, indicating that  $CPU_j$  is idle. The arcs going from transitions  $t_{1,j}^{exec}$ , ...,  $t_{n,j}^{exec}$  to place  $p_j^{idle}$  and from place  $p_j^{idle}$ to transitions  $t_{1,j}^{alloc}$ , ...,  $t_{n,j}^{alloc}$  are weighted by a constant value  $\eta$ , to ensure that the flow in transitions  $t_{1,j}^{alloc}$ , ...,  $t_{n,j}^{alloc}$  is limited by the throughput CPU capacity modeled by place  $p_j^{idle}$ .

Thermal model.- The execution of cycles of  $\tau_i$  in  $CPU_j$  is modeled by firing transitions  $t_{i,j}^{exec}$  (Fig. 1), which adds tokens to places  $p_i^{com_j}$ , activating the thermal model for  $CPU_j$  (whose temperature will increase because of the activity). The rest of transitions and places in this thermal model represent heat transport by conduction and convection, for a deeper explanation see Desirena-Lopez et al. (2014).

Task and thermal model evolution.- The dynamic behavior of the global model (Fig. 1) is provided by the following equations:

$$\dot{m}_T = C_T \Lambda_T \Pi_T(m) m_T + C_a \Lambda_a \Pi_a(m) m_a + C^{exec}_{exec} f^{exec}$$

$$\dot{\boldsymbol{n}}_{\boldsymbol{a}} = 0$$
 (3b)

(3a)

$$\dot{m}_{\mathcal{T}} = C_{\mathcal{T}} \Lambda_{\mathcal{T}} \Pi_{\mathcal{T}}(m) m_{\mathcal{T}} - C_{\mathcal{T}}^{alloc} w^{alloc}$$
 (3c)

$$\dot{m}_{\mathcal{P}} = C_{\mathcal{P}} \Lambda_{\mathcal{P}} \Pi_{\mathcal{P}}(m) m_{\mathcal{P}} + C_{\mathcal{P}}^{alloc} w^{alloc}$$
(3d)

$$h_{exec} = C_{\mathcal{P}}^{exec} f^{exec}$$
(3e)

 $C_x$ ,  $\Lambda_x$ , and  $\Pi_x(m)$  are the incidence matrix, the firing rate transitions and the configuration matrix respectively

ń

 $(x = \{T, a, \mathcal{T}, \mathcal{P}\})$  of the thermal, task and processors subnet.  $C_T^{alloc}$ ,  $C_T^{alloc}$  and  $C_P^{alloc}$  stand for the connections of transitions  $t_{i,j}^{alloc}$  from (to) places in the thermal model, task and CPU throughput model, respectively. Matrix  $C_{\mathcal{P}}^{exec}$  holds the columns of the transitions  $t_{i,j}^{exec}$  of the incidence matrix  $C_{\mathcal{P}}.w^{alloc}$  is the controlled flow of the *allocation transitions* (i.e. the task allocation rate to CPUs). Eq. (3a) represents the system's temperature evolution. Eq. (3b) indicates that the environmental temperature keeps constant at all times (its derivative is neglected). Eq. (3c) describes the arrival of periodic tasks to the system. Eq. (3d) models the CPUs cycles that are assigned to tasks. Finally, Eq. (3e) models task execution.

#### 5. OFF-LINE STAGE

In our approach, prior to considering the aperiodic tasks, we need to find the minimum and maximum frequencies to execute the periodic task set subject to temporal and thermal constraints. This requires first to study the system thermal behavior. Later, we will compute a schedule applying a deadline partitioning approach.

Steady state Thermal Analysis.- Task execution generates a thermal activity given by Eq. (3a), where  $C_P^{exec} = F^3 C_P^{(exec)}$ . This can be rewritten as a space state equation:  $\dot{m}_T = Am_T + B'm_a + F^3 B f^{exec}$  and  $Y_T = Sm_T$ (the CPUs temperature), where  $A = C_T \Lambda_T \Pi_T(m)$ ,  $B = C_P^{(exec)}$  and  $B' = C_a \Lambda_a \Pi_a(m)$ . Since the schedule is periodic, the temperature is a non-decreasing function reaching a steady state temperature  $(m_{T_{ss}})$ , i.e.  $\dot{m}_T = 0$  when time tends to infinite. Hence  $m_{T_{ss}} = -A^{-1}(F^3 B w^{alloc} + B'm_a)$ .

Since  $Sm_{T_{ss}} \leq T_{max}$  then:

$$-SA^{-1}F^{3}Bw^{alloc} \le T_{max} + SA^{-1}B'm_{a}$$

$$\tag{4}$$

This equation provides the thermal constraints that the allocation of tasks to the processors  $(\boldsymbol{w^{alloc}})$  must fulfill.

Energy consumed by a schedule.- If the  $CPU_j$  clock frequency is F during the time interval  $(\zeta_1, \zeta_2]$ , then the average energy consumed during this interval by the tasks running on  $CPU_j$  is defined as:

$$E_j = \int_{\zeta_1}^{\zeta_2} P_{CPU_j}(F) d_{\zeta} \tag{5}$$

 $P_{CPU_j}(F)$  is the power consumed by a  $CPU_j$ . It depends on  $P_{dyn}(F)$ , the dynamic power due to computational activities of tasks, and  $P_{leak}$ , the static power due to leakage. It is computed as:  $P_{CPU_j}(F) = P_{dyn_j}(F) +$  $P_{leak_j} = \lambda_{i,j}^{exec} m_{i,j}^{busy}(\zeta)F^3 + P_{leak_j}$ . Where  $P_{leak_j}$  can be modeled as a linear function of temperature (Ahmed et al. (2016)):  $P_{leak} = \delta T + \rho$ , where T is the CPUs temperature and  $\delta$  and  $\rho$  are modeling constants. The consumed energy is minimized *iff* the clock frequency F is minimized, but F must be high enough to ensure that temporal constraints are met. Next, we compute this frequency. Minimum frequency Modern CPUs vary their clock frequency according to a number of preset values, i.e.  $\mathcal{F} = [F_1, ..., F_{max}]$ . We normalize this set as  $\phi = [\phi_{min} = \frac{F_1}{F_{max}}, ..., 1]$ . The next proposition obtains the minimum clock frequency that fulfills temporal constraints.

*Proposition 5.1.* Assuming that the task utilization is less than the number of processors in the METARTS problem, the normalized clock frequency that minimizes the total energy consumption while meeting temporal constraints is constant:

$$\Phi^* = \max\{\phi_{min}, \frac{1}{m} \sum_{i=1}^n \frac{cc_i}{\omega_i F_{max}}\}$$
(6)

Proof 5.1. According to Eq. (5), the energy has a minimum iff the consumer power is minimum. This occurs when  $\phi^3$  is minimum and fulfills that  $\sum_{i=1}^{n} \frac{cc_i}{\phi w_i F_{max}} = m$ , and  $\phi \ge \phi_{min}$ . Using Lagrange multipliers, the Lagrangian function is  $L = \phi^3 + \mu_1(\frac{1}{\phi}\sum_{i=1}^{n}\frac{cc_i}{w_i F_{max}} - m) + \mu_2(\phi_{min} - \phi)$ . The solution yields four cases: a) Both multipliers are inactive  $(\mu_{1,2} = 0)$ ; b) Both multipliers are active  $(\mu_{1,2} \ge 0)$ ; c)  $\mu_1 = 0$  and  $\mu_2 \ge 0$ ; and d)  $\mu_1 \ge 0$  and  $\mu_2 = 0$ . The first case is unfeasible, because  $\phi$  cannot be zero. In the second case, the only solution is  $\phi = \phi_{min} = \frac{1}{m} \sum_{i=1}^{n} \frac{cc_i}{w_i F_{max}}$ . Finally, if one multiplier is active while the other one is inactive there are two possible solutions:  $\phi = \phi_{min}$  or  $\phi = \frac{1}{m} \sum_{i=1}^{n} \frac{cc_i}{w_i F_{max}}$ . Consequently, in order to fulfill both constraints, the normalized clock frequency that minimize the total energy consumption becomes  $\Phi^* = \max\{\phi_{min}, \frac{1}{m}\sum_{i=1}^{n} \frac{cc_i}{\omega_i F_{max}}\}$ .  $\Box$ 

The normalized frequency  $\Phi^*$  meets the temporal constraints. To guarantee that the thermal constraints are also fulfilled, we must compute  $\boldsymbol{w}^{alloc}$  and solve Eq. (4). At frequency  $\Phi^*$  the total system utilization is  $U = \sum_{i=1}^{n} \frac{cc_i}{\omega_i \phi^* F_{max}} = m$  and the processor frequency is  $F^* = \min\{F \in \mathcal{F} | F \ge \Phi^* F_{max}\}$ , given the nature of the discrete set of frequencies. Since we have a fully utilized system, the distribution of the CPU cycles required to execute all tasks must be homogeneous, i.e.,  $\boldsymbol{w}^{alloc} = [\frac{1}{m} \sum_{i=1}^{n} \frac{cc_i}{\omega_i F^*}, \ldots, \frac{1}{m} \sum_{i=1}^{n} \frac{cc_i}{\omega_i F^*}]^T$ . Moreover, if  $\boldsymbol{w}^{alloc}$ satisfies Eq. (4), then the thermal constraints are also satisfied. Otherwise, the METARTS problem does not have a solution. If it has a solution ( $\Phi^*$  is feasible), then we can compute the maximum CPU cycles available for aperiodic tasks, and the maximum clock frequency that can be used subject to thermal constraints.

Maximum CPU cycles and clock frequency The maximum thermal frequency  $F^+$  is the greatest frequency at which all CPUs operate at 100% of utilization and the temperature never exceeds the maximum thermal constraints.  $F^+$  can be computed by using the next programming problem. The first constraint is thermal.  $CC_j$  represents the cycles that  $CPU_j$  must execute per hyperperiod. Since all CPUs must work at their maximum capacity, the second constraint implies that the CPU utilization is 100%. The last constraint bounds  $F^+$  to the actual clock frequency range in the MPSoC.  $\begin{array}{ll}max & F^+\\s.t.\end{array}$ 

$$-SA^{-1}F^{+3}B\left[\frac{CC_1}{F^+H}, \dots, \frac{CC_m}{F^+H}\right]^T \leq T_{max} + SA^{-1}B'm_a$$
$$\frac{CC_j}{F^+H} = 1 \qquad \forall j = 1, \dots, m$$
$$F^* \leq F^+ \leq F_{max}$$
(7)

The solution for  $F^+$  has to be in the set  $\mathcal{F}$  of discrete frequencies. Thus the processor frequency is updated as  $F^+ = \max\{F \in \mathcal{F} | F \leq F^+\}.$ 

#### 5.1 Deadline partitioning

We consider the ordered set of all tasks' jobs deadlines to define scheduling intervals, as in deadline partitioning (Funk et al. (2011)). Each task  $\tau_i$  must be executed  $n_i = \frac{H}{\omega_i}$  times within the hyperperiod H. Thus every  $q * \omega_i$ , where  $q = 1, ..., n_i$  is a deadline that must be considered in the analysis. These deadlines can be ordered and joined in the set  $SD_i = \{sd_i^1, ..., sd_i^{n_i}\}$ . A general set of deadlines is defined as  $SD = SD_0 \cup ... \cup SD_{|\mathcal{T}|}$ where  $SD_0 = \{0\}$ . The elements of SD can be arranged in ascendant order and renamed as  $SD = \{sd_0, ..., sd_\alpha\},\$ where  $\alpha$  is the last deadline. The scheduling interval  $I_{SD}^k =$  $[sd_{k-1}, sd_k]$  is defined and  $|I_{SD}^k| = sd_k - sd_{k-1}$  represents the scheduling interval duration. The proposed deadline partitioning problems assumes a 100% utilization on every CPU, but recalling that  $F^* = \min\{F \in \mathcal{F} | F \ge \Phi^* F_{max}\},\$ in most cases  $F^* \neq \Phi^* F_{max}$  consequently idle cycles appear. To solve this problem we introduce a dummy task with period H and utilization  $m - \sum_{i=1}^{n} \frac{c_i}{F^* \omega_i}$ . Then the cycles that each task must execute in the  $I_{SD}^k$ , i.e  $x_i^k$ , can be computed on fall. be computed as follows.

Let  $cc_i^* = \omega_i * F^* - cc_i$  be the cycles that task  $\tau_i$  can be idle. Thus, the total amount of cycles  $(sd_k * F^*)$  in  $sd_k$ can be rewritten as  $sd_k * F^* = q * \omega_i * F^* + r_i$ , where  $0 \leq r_i < \omega_i * F^*$  and  $q \in \mathbb{Z}$ , where q represents the occurrences of a task (the task's jobs) in the system. If  $r_i = 0$ , it means that  $\tau_i$  has its deadline in the scheduling interval. Then the following LPP can be posed to compute  $x_i^k$ .

$$\min \sum_{i=1}^{n} x_{i}^{k} \quad \forall k = 1, \dots, \alpha$$

$$s.t$$

$$\forall k \sum_{i=1}^{n} x_{i}^{k} = m * |I_{SD}^{k}| * F^{*}$$

$$\text{if } r_{i} = 0 \sum_{\gamma=1}^{k} x_{i}^{\gamma} = q * cc_{i}$$

$$\text{if } r_{i} \neq 0 \sum_{\gamma=1}^{k} x_{i}^{\gamma} \ge -q * cc_{i}^{*} + max\{0, \sum_{\gamma=1}^{k} |I_{SD}^{\gamma}| * F^{*} - cc_{i}^{*}\}$$

$$\forall i \quad x_{i}^{k} \le |I_{SD}^{\gamma}| * F^{*}$$

$$(8)$$

The first constraint implies that the CPU utilization is 100%. It is required since  $\Phi^*$  indicates that CPU utilization is 100%. The second constraint guarantees that

those tasks that must complete execution in this interval actually end. The last constraint guarantees deadline fulfillment. The following proposition guarantees that if the former LLPs are orderly solved according to the k - thinterval, then the computed amount of time that each task must run per interval yields a feasible schedule.

Proposition 5.2. Given a task set  $\mathcal{T}$  presented in Definition 3.1, where the task utilization at  $F^*$  is equal to the number of CPUs, the solution of the linear programming problems in Eq. (8) is always integer and if each task  $\tau_i$  is executed exactly  $x_i^k$  cycles during the k-th interval, then a feasible schedule is obtained.

Proof 5.2. Let  $\mathcal{T}^k = \mathcal{T}^k_1 \cup \mathcal{T}^k_2$  be a task set partition, where  $\mathcal{T}^k_1 = \{\tau_1, ..., \tau_v\}$  is the set of tasks that have their deadlines at  $sd_k$  and  $\mathcal{T}^k_2 = \{\tau_{v+1}, ..., \tau_n\} = \mathcal{T} - \mathcal{T}^k_1$ . Some slack variables  $h_i$  are added to form a set of constraints that can be represented as My = b where  $y = [x h]^T$  (i.e., equality constraints are obtained). Notice that vector bis always integer. By construction, the restriction matrix

*M* has the form:  $M = \left[\frac{L_{(v+1)\times n}}{Q_{(2n-v)\times n}} | \frac{\varnothing}{I_{(2n-v)}}\right]$ , and L has the form:  $L_{(v+1)\times n} = \left[\frac{1 \cdots 1}{I_{v\times v} | \varnothing}\right]$ . It is easily seen that

rank(L) = v + 1 and rank(M) = rank(L) + rank(I) =2n + 1, i.e M is a full row rank. The solution is always integer if M is unimodular, i.e., the determinant of every square submatrix  $(M_{si})$  of M is either 0,+1 or -1. All  $M_{si}$ are obtained deleting columns, and there are three possible scenarios: first, if any of the first v columns is removed,  $M_{si}$  loses rank, hence  $det(M_{si}) = 0$ . Second, if any deleted column contains a nonzero entrie where its corresponding row has a nonzero element among the first v columns,  $M_{si}$ loses rank since the resulting row is duplicated among the first v rows. Thus,  $det(M_{si}) = 0$ . Finally, when any other column not listed before is deleted, the resulting matrix always can be arranged as  $M_{si} = \begin{bmatrix} A | \varnothing \\ B | I \end{bmatrix}$ , where matrix A is always TUM, according to **Theorem 3.4** reported in Sierksma (2001), thus  $det(A) = 0, \pm 1$ . Therefore,  $det(M_{si}) = det(A) \cdot det(I) = 0, \pm 1. \quad \Box$ 

## 6. ON-LINE STAGE: SCHEDULER

The previous section computed the CPU clock frequency  $F^*$ , maximum clock frequency  $F^+$  and task execution time per scheduling interval. The on-line scheduler uses these data to implement a Fixed Priority Zero-Laxity (FPZL) algorithm (Davis and Burns (2011a)). It allocates tasks' jobs during their respective scheduling interval upon the occurrence of three possible events: a zero-laxity event (a job must immediately execute lest it misses its deadline), job completion or the arrival of an aperiodic task. During the  $I_{SD}^k$  interval task  $\tau_i^k$  must execute  $x_i^k$  cycles at a given  $F_n$  clock frequency.

Whenever an event occurs, the task pri-Priority Levels ority is updated as follows. Jobs reaching their zero-laxity time are given the maximum priority (= 1). Jobs being executed and with laxity different from zero receive priority equal to 2. The remaining jobs receive priority level equal to 3 (the lowest one). Thus, zero laxity tasks have the highest priority and must be executed immediately.

## Algorithm 1 On-line schedule

- 1: Input  $I_{SD}^k$  Schedulinng intervals;  $X^k$  tasks CPU cycles per interval;  $ex_i^k$  execution CPU cycles per interval
- Output A feasible schedule;
  - $k = 0, \zeta = 0$
- 3: Compute the ordered set of laxities as:  $SL = \{l_i | l_i = sd_{k+1} (F_n * x_i^k ex_i^k) \zeta\}$

$$SL = \{l_i | l_i = sd_i\}$$

- 4: while  $\zeta \leq H$  do 5:
- while  $\zeta \leq sd_{k+1}$  do 6: Compute task priorities using Priority Levels
- 7: Execute the m tasks with higher priority until an event occurs (An event occurs if a task reaches its zero laxity, task ends or aperiodic task arrives.)
- Compute the ordered set of laxities as: 8:  $SL = \{l_i | l_i = sd_{k+1} - (F_n * x_i^k - exec_i^k) - \zeta\}$  $\zeta = \zeta + \text{ current time}$
- 9:
- 10: end while

11. k = k + 1

Execution of m tasks with the highest priority In Alg. 1 step 8, *m* tasks must be executed (i.e. allocated to a CPU). In order to reduce the number of migrations, tasks that are executed during two consecutive events are allocated to the same CPU.

## 7. APERIODIC TASKS

Aperiodic tasks arrive asynchronously to the system. The system determines if these tasks can be executed without compromising the hard real-time constraints of the periodic task set. If so, a new CPU clock frequency is computed to allow the execution of the aperiodic task. This frequency must be in the range  $\mathcal{F}^s = \{F^* \dots F^+\}$  (every frequency in this range meets the thermal constraints), and it is kept as low as possible to guarantee a minimum energy consumption while meeting the temporal constraints.

Fig. 2 shows the scheme to schedule aperiodic tasks without compromising the hard real-time constraints of the periodic task set. The off-line stage computes the optimal and maximum allowed frequencies  $F^*$ ,  $F^+$ , the scheduling intervals  $I_{SD}^k$  and the periodic task CPU cycles that must be executed per scheduling interval  $X = \{X^1, ..., X^{\alpha}\},\$ where  $X^k = \{x_1^k, \dots, x_n^k\}$  represents the set of *CPU* cycles  $x_i^k$  that must be executed during the k - th scheduling interval of task  $\tau_i$ . Job execution is tracked during the  $k\!-\!th$  interval by the system and passed as the input to the adaptive scheduler (AS). AS is activated at the arrival or ending time of an aperiodic task  $\tau_i^a$ . At this time, the task's CPU cycles  $cc_i^a$  and its relative deadline  $d_i^a$  are sent to AS. This information together with the outputs of the offline stage are used to compute the new frequency and the *CPU* cycles  $(x_{\tau_i^a}^k)$  required to execute  $\tau_i^a$  per scheduling interval, such that  $X^k = X^k \cup \{x_{\tau_i^a}^k\}$ . When  $\tau_i^a$  finishes its execution, the frequency is recalculated and the CPUcycles demanded by  $\tau_i^a$  are discarded.

The complexity of the On-line stage depends Complexity on two algorithms. The priority level and the computation of laxity in Alg. 1 is linear in the number of tasks. At most  $n = |\mathcal{T}|$  tasks will end its execution  $x_i^k$  in the k-th interval (there at most n tasks). Also at most n tasks will reach their zero laxity. If q aperiodic tasks arrive in the k-thinterval, then the nested while loop ends in  $(n + n + q) \times$ 



Fig. 2. Minimum Energy Thermal Aware scheme for periodic and aperiodic tasks.

(n+n) (number of events  $\times$  number of operations). Considering that the outer loop runs  $\alpha = |I_{SD}|$  times, then the number of steps of this algorithm is polynomial in the order of tasks. Alg. 2 runs on the arrival of an aperiodic task and is polynomial in the order of tasks and independent of the number of CPUs. Thus the proposed algorithm is polynomial in the order of tasks.

Algorithm 2 Adaptive Scheduler (AS)

- 1: Input  $cc_i^a, d_i^a$  Aperiodic tasks parameters;  $ex_i^k$  Cycles executed in the system for all active tasks.
- **Output** New Frequency  $F_n \in \mathcal{F}^s = \{F^*, \ldots, F^+\}$ , task *CPU* cycles per  $I^k_{SD}, x^k_{\tau_{a_i}}$  Per-interval CPU cycles for execution of the aperiodic task.
- **Initialize**  $n = |\mathcal{T}|, m = |\mathcal{P}|, F_n = F^*, q$  the aperiodic tasks 3: currently being attended.
- BEGIN
- 4: if aperiodic task arrives then
- 5:Let  $r_i^a = \text{current time when } \tau_i^a \text{ arrives};$ Compute required *CPU* cycles for active tasks from  $r_i^a$  to 6:

$$\begin{split} r_i^a + d_i^a; \\ C_u &= \sum_{i=1}^{|X^k|} (x_i^k - e x_i^k) + \sum_{\gamma=k+1}^{\Gamma} \sum_{i=1}^{|X^{\gamma}|} x_i^{\gamma}; \end{split}$$

where k is the current scheduling interval at  $r_i^a$  and and  $\Gamma$  is the scheduling interval at  $r_i^a + d_i^a$ ;

- $C_{free} = m * d_i^a * F^+ C_u$ ; the free CPU cycles in the interval 7:  $[r_i^a, r_i^a + d_i^a];$
- $\mathbf{if} C_{free} \geq c c_i^a \mathbf{then}$ 8:
- Accept task  $\tau_i^a$ ; 9:

10:  

$$F_{n} = \min F \in \mathcal{F}^{s} \text{ such that } F \geq \frac{C_{u} + cc_{i}^{a}}{m * d_{i}^{a}};$$

$$cc_{r} = cc_{i}^{a};$$

$$x_{\tau_{i}}^{k} = \min \left\{ m(|I_{SD}^{k}| - r_{i}^{a})F_{n} - \sum_{i=1}^{|X^{k}|} (x_{i}^{k} - ex_{i}^{k}), cc_{r} \right\};$$
in the  $k - th$  interval;  
for  $\gamma = k + 1$  to  $\Gamma$  do - in other intervals  

$$cc_{r} = cc_{r} - x_{\tau_{i}}^{\gamma-1};$$

$$x_{\tau_{i}}^{\gamma} = \min \left\{ m(|I_{SD}^{\gamma}|)F_{n} - \sum_{i=1}^{|X^{\gamma}|} x_{i}^{\gamma}, cc_{r} \right\};$$
end for  
11: else  
12: Reject task;  
13: end if  
14: end if  
15: if an aperiodic task finishes then

Discard the CPU cycles associated to the aperodic task Recalculate the new frequency;

end if END

## 8. SIMULATIONS RESULTS

In order to show how to use the proposed scheme, a proof of concept is presented. It consists of a set of sporadic periodic tasks  $\mathcal{T} = \{\tau_1, \tau_2, \tau_3\}$ , where  $\tau_1 = (2000, 4), \tau_2 =$  $(5000, 8), \tau_3 = (6000, 12),$  the hyperperiod is H = 24. These tasks run on two homogeneous microprocessors where the isotropic thermal properties and dimensions of the materials are taken from Desirena-Lopez et al. (2014). The processor supports four operating frequency levels  $\mathcal{F} = \{0.5, 0.85, 0.95, 1\} KHz$ . The temperature of the surrounding air is set to  $45^{\circ}C$  and it is constant. The maximum operating temperature level is set to  $T_{max_{1,2}} =$  $50^{\circ}C$ . The simulations herein presented consider CPUs with caches and speculative mechanisms non-existent or turned off. First, the minimum frequency for the periodic task set is off-line computed according to Eq. (6), obtaining  $\Phi^*~=~0.8125,$  hence the selected frequency is  $F^*~=$ 0.85 kHz. Eq. (7) provides the maximum clock frequency  $(F^+ = 1kHz)$ , so that the METARTS problem has a solution. We assume that sheeduling and context switch overheads are included in the tasks' WCET. Then, solving the LPP in Eq. (8) for  $F^*$  yields the CPU cycles of each task to be executed at each interval  $(x_i^k)$ . Fig. 3 provides the schedule and temperature evolution produced by the algorithm without considering aperiodic tasks. Fig. 4 depicts the outcome of the algorithm and the evolution of the temperature when an aperiodic task  $\tau_1^a =$ (2000, 10) arrives at  $\zeta = 2$ , during the  $I_{SD}^{1}$  interval, thus  $\tau_{1}^{a}$  has an absolute deadline at  $\zeta = 12$ . Since  $C_{free} \geq cc_{1}^{a}$ , the AS accepts the aperiodic task, and computes  $F_{n} =$  $950kHz \in \mathcal{F}^s$  as the frequency at which the processors must execute during interval [2, 12]. Fig. 4 shows that temperature increases during this interval because of the execution of the tasks at frequency  $F_n$ , and then it decreases after  $\zeta = 12$  because a new (lower) frequency has been calculated for the next interval. In both experiments, with and without the aperiodic tasks,  $CPU_1$  achieves full utilization, whereas  $CPU_2$  shows a slack (idle time) at about  $\zeta = 18$ , which translates into a temperature valley. This slack appears because the exact optimal frequency calculated in Eq. (6) is ceiled to a frequency belonging the discrete set of frequencies available in the microprocessor  $(F_n \in \mathcal{F}^s).$ 

#### 9. CONCLUSIONS

This work shows that the TCPN formalism is a suitable way to model real-time task scheduling problems considering thermal, temporal and energy restrictions. Upon a TCPN model, we build a two-stage thermal-aware realtime system in which a periodic task set executes at minimum clock frequency to save energy and achieve maximum CPU utilization while honoring thermal constraints. Aperiodic tasks are also dynamically accepted as long as increasing the clock frequency allows to obtain a suitable slack. Time and thermal constraints are preserved in all cases. We have assumed that the aperiodic task  $(\tau_i^a)$ fits in the slack inside the hyperperiod  $(r_i^a + d_i^a \leq H)$ . Immediate further work includes to relax that condition, and to improve our slack reclaiming approach to decrease context switching when accepting aperiodic tasks. Also, the fact that our underlying TCPN model allows modeling



Fig. 3. Temperature evolution (upper plot) for the periodic schedule (lower plot) at  $CPU_1$  (above) and  $CPU_2$ (below). The maximum temperature produced by this schedule is  $T_{CPU_{1,2}} = 46.5^{\circ}C$ 



Fig. 4. Temperature evolution (upper plot) for the periodic schedule (lower plot) at  $CPU_1$  (above) and  $CPU_2$ (below) upon acceptance of the aperiodic task  $\tau_1^a$ . The maximum temperature produced by this schedule is  $T_{CPU_{1,2}} = 47.06^{\circ}C$ 

resource sharing, a tough problem in multicore real-time scheduling, opens up a promising venue. Last, we still have to measure how much energy we can save with respect to other scheduling techniques for aperiodic tasks.

## ACKNOWLEDGEMENTS

This work was partially supported by grants TIN2016-76635-C2-1-R (AEI/FEDER, UE), gaZ: T48 research group (Aragón Gov. and European ESF), and HiPEAC4 (European H2020/687698)

# REFERENCES

- Ahmed, R., Ramanathan, P., and Saluja, K.K. (2016). Necessary and sufficient conditions for thermal schedulability of periodic real-time tasks under fluid scheduling model. ACM Transactions on Embedded Computing Systems (TECS), 15(3), 49.
- Baruah, S., Bertogna, M., and Butazzo, G. (2015). Multiprocessor Scheduling for Real-Time Systems. Springer-Verlag New York, Inc., Secaucus, NJ, USA.

- Baruah, S.K., Cohen, N.K., Plaxton, C.G., and Varvel, D.A. (1996). Proportionate progress: A notion of fairness in resource allocation. *Algorithmica*, 15(6), 600– 625.
- David, R. and Alla, H. (2008). Discrete, continuous and hybrid Petri nets (david, r. and alla, h.; 2004). Control Systems, IEEE, 28(3), 81–84.
- Davis, R.I. and Burns, A. (2011a). Fpzl schedulability analysis. In Proceedings of the 2011 17th IEEE Real-Time and Embedded Technology and Applications Symposium, RTAS '11, 245–256. IEEE Computer Society, Washington, DC, USA.
- Davis, R.I. and Burns, A. (2011b). A survey of hard real-time scheduling for multiprocessor systems. ACM computing surveys (CSUR), 43(4), 35.
- Desirena-Lopez, G., Briz, J.L., Vázquez, C.R., Ramírez-Treviño, A., and Gómez-Gutiérrez, D. (2016). On-line scheduling in multiprocessor systems based on continuous control using timed continuous petri nets. In 13th International Workshop on Discrete Event Systems, 278–283.
- Desirena-Lopez, G., Vázquez, C.R., Ramírez-Treviño, A., and Gómez-Gutiérrez, D. (2014). Thermal modelling for temperature control in MPSoC's using fluid Petri nets. In *IEEE Conference on Control Applications part* of Multi-conference on Systems and Control.
- Funk, S., Levin, G., Sadowski, C., Pye, I., and Brandt, S. (2011). Dp-fair: a unifying theory for optimal hard realtime multiprocessor scheduling. *Real-Time Systems*, 47(5), 389–429.
- Hettiarachchi, P.M., Fisher, N., Ahmed, M., Wang, L.Y., Wang, S., and Shi, W. (2014). A design and analysis framework for thermal-resilent hard real-time systems. *Embedded Computing Systems, ACM Transactions on*, 13(5s), 146:1–146:25.
- Kong, J., Chung, S.W., and Skadron, K. (2014). Recent thermal management techniques for microprocessors. *ACM Computing Surveys*, 44(3), 13:1–13:42.
- Schor, L., Bacivarov, I., Yang, H., and Thiele, L. (2012). Worst-case temperature guarantees for real-time applications on multi-core systems. In 2012 IEEE 18th Real Time and Embedded Technology and Applications Symposium, 87–96. IEEE.
- Sierksma, G. (2001). Linear and integer programming: theory and practice. CRC Press.
- Silva, M., Júlvez, J., Mahulea, C., and Vázquez, C.R. (2011). On fluidization of discrete event models: observation and control of continuous Petri nets. *Discrete Event Dynamic Systems*, 21(4)(3), 427–497.
- Silva, M. and Recalde, L. (2007). Redes de Petri continuas: Expresividad, análisis y control de una clase de sistemas lineales conmutados. *Revista Iberoamericana* de Automática e informática Industrial, 4(3), 5–33.