

# ERPOT: A Quad-Criteria Scheduling Heuristic to Optimize Execution Time, Reliability, Power Consumption and Temperature in Multicores

Athena Abdi, Alain Girault, Hamid Zarandi

## ► To cite this version:

Athena Abdi, Alain Girault, Hamid Zarandi. ERPOT: A Quad-Criteria Scheduling Heuristic to Optimize Execution Time, Reliability, Power Consumption and Temperature in Multicores. IEEE Transactions on Parallel and Distributed Systems, Institute of Electrical and Electronics Engineers, 2019, 30 (10), pp.2193-2210. 10.1109/TPDS.2019.2906172 . hal-02400019

**HAL Id: hal-02400019**

**<https://hal.inria.fr/hal-02400019>**

Submitted on 9 Dec 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# ERPOT: A Quad-Criteria Scheduling Heuristic to Optimize Execution Time, Reliability, Power Consumption and Temperature in Multicores

Athena Abdi \* Alain Girault ‡ Hamid Zarandi \*

**Abstract**—We investigate multi-criteria optimization and Pareto front generation. Given an application modeled as a Directed Acyclic Graph (DAG) of tasks and a multicore architecture, we produce a set of non-dominated (in the Pareto sense) static schedules of this DAG onto this multicore. The criteria we address are the execution time, reliability, power consumption, and peak temperature. These criteria exhibit complex antagonistic relations, which make the problem challenging. For instance, improving the reliability requires adding some redundancy in the schedule, which penalizes the execution time. To produce Pareto fronts in this 4-dimension space, we transform three of the four criteria into constraints (the reliability, the power consumption, and the peak temperature), and we minimize the fourth one (the execution time of the schedule) under these three constraints. By varying the thresholds used for the three constraints, we are able to produce a Pareto front of non-dominated solutions. We propose two algorithms to compute static schedules. The first is a ready list scheduling heuristic called ERPOT (Execution time, Reliability, POver consumption and Temperature). ERPOT actively replicates the tasks to increase the reliability, uses Dynamic Voltage and Frequency Scaling to decrease the power consumption, and inserts cooling times to control the peak temperature. The second algorithm uses an Integer Linear Programming (ILP) program to compute an optimal schedule. However, because our multi-criteria scheduling problem is NP-complete, the ILP algorithm is limited to very small problem instances. Comparisons showed that the schedules produced by ERPOT are on average only 10% worse than the optimal schedules computed by the ILP program, and that ERPOT outperforms the PowerPerf-PET heuristic from the literature on average by 33%.

**Keywords**—multicore static scheduling, Reliability, Temperature, Power consumption, Multi-objective optimization, Pareto front.

## I. INTRODUCTION

Multicores are widely used in modern safety critical embedded systems design. Their advantages over super-scalar processor architectures are lower power consumption, higher performance, and lower design complexity [1]. When designing safety critical applications, many non-functional criteria must be addressed. The most important ones are the *total execution time* (because these systems must react to inputs within a fixed delay), the *reliability* (because failures could have fatal consequences), the *power consumption* (to maximize the autonomy of the system when it operates on a battery),

and the *temperature* (because of its negative influence on processing speed, reliability, and power consumption) [1]–[4]. There are many real-life applications that motivate our study, including satellite systems, portable medical devices, and full authority digital engine control (FADEC) in aircraft.

Considering these four criteria simultaneously during the design phase is very difficult because they are *antagonistic* [1], [2], [4]–[9]. For instance, the total execution time and reliability are antagonistic because increasing the reliability requires some form of redundancy (be it spatial or temporal), which negatively impacts the execution time. Similarly, the execution time and the temperature are antagonistic because adding idle times to cool the cores obviously has a negative impact on the execution time. Finally, the execution time and the power consumption are antagonistic because reducing the power consumption requires lowering the operating voltage and frequency of the cores, which increases the execution time. Those tradeoffs are easy to grasp (but difficult to address), but other tradeoffs are less obvious: for instance, lowering the operating voltage and frequency of a core (which lowers the power consumption) increases the nominal failure rate per time unit of this core. The reason is that the sensitivity of processors to energy particles leads to an increase of the failure rate at low voltage/frequency operating points [10], [11], because lowering the voltage decreases the critical charge of the circuit. As a consequence, the power consumption and the reliability are also antagonistic. Failing to take into account these antagonisms could result in bad design choices.

These antagonisms call for the computation of as many tradeoffs as possible, rather than a single tradeoff, so that the user will have a choice. We must therefore produce a *set* of solutions in the 4-dimensions space (execution time, reliability, power consumption, temperature). We rely on the notion of *Pareto dominance*, and we use a variant of the  $\epsilon$ -constraint method [12], [13] coupled with a scheduling algorithm that accounts for the four criteria to produce the *Pareto front* in this 4D space. More precisely, we transform three criteria into *constraints* (the reliability, the power consumption, and the peak temperature), and we *minimize* the fourth one (the execution time of the schedule) under these three constraints.

Although several studies have addressed some of these parameters, none have considered these four criteria jointly in an optimization problem. For instance, some studies completely ignore the reliability [4], [14] or the temperature [2], [9]. Other studies tackle the problem as a hardware/software co-design problem, jointly optimizing the floorplan of the multicore and

\* Department of Computer Engineering and Information Technology, Amirkabir University of Technology (Tehran Polytechnic), Tehran, Iran.

‡ Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LIG, 38000 Grenoble, France.

the schedule of the application task graph to minimize the peak temperature [14], but without considering the reliability.

We therefore propose a static scheduling heuristic method called ERPOT, an acronym that stands for Execution time, Reliability, POver consumption and Temperature. Given an application modeled as a Directed Acyclic Graph (DAG) of tasks, a multicore architecture, and thresholds on the reliability, the power consumption, and the temperature, ERPOT generates a static schedule of this DAG onto this multicore such that each constraint is below its corresponding threshold, and such that the execution time is as small as possible. Each schedule is interpreted as a point in the 4D space (execution time, reliability, power consumption, temperature). By varying the values of the thresholds and calling iteratively ERPOT, we are able to produce a full Pareto front in this 4D space.

The problem of scheduling a DAG of tasks onto a distributed architecture is known to be NP-complete [15], and so is the multi-criteria scheduling problem, which motivates the design of a heuristic algorithm. Additionally, we present an ILP program of the optimization problem, which is used to validate ERPOT (i.e., both algorithms produce the same schedule on the same problem instance) and to assess experimentally how good ERPOT is. Comparing the results of ERPOT with the optimal results obtained by the ILP program shows that the average difference is less than 10%. However, ERPOT is much faster than the ILP program, which fails to complete even for application graphs of relatively small sizes (8 tasks at most).

The key contributions of this paper are:

- The ERPOT quad criteria scheduling heuristic, which optimizes the *execution time*, the *reliability*, the *power consumption*, and the *temperature*.
- A 4D variant of the  $\varepsilon$ -constraint method [12] to build the Pareto front of the solutions in the 4D space (execution time, reliability, power, temperature).
- An ILP program of the quad criteria optimization problem to compare the solution computed by ERPOT with the optimal solution.

ERPOT extends the heuristics proposed in [2] by taking into account the peak temperature. The first challenge of doing so lies in the intricate dependence of the temperature on the other criteria of [2], namely the failure rate, the power consumption, and the execution time. The second challenge is in the scheduling heuristic itself: each scheduling decision is made by “predicting” what will be the value of the temperature, power consumption, and failure rate at the end of the task being scheduled. However, the temperature *varies* during the execution of the task, because it obeys the classical thermal differential equation. Since the power consumption (and similarly the failure rate) depends on the temperature, the computation of the power consumption is inexact unless it is performed continuously during the execution of the task being scheduled, which is much too expensive. Addressing this challenge requires an over-approximation of the temperature and the proof that this is safe for the power-consumption constraint. This was not the case when only the power consumption and the failure rate were considered, making the scheduling heuristic of [2] much simpler. The third challenge resides in maintaining the peak temperature below a given threshold, which involves

a combination of lowering the voltage/frequency (thanks to DVFS), inserting cooling intervals, and over-estimating the temperature when there are “holes” at the end of schedule under construction. A final contribution compared to [2] is that The ILP program of [2] does not consider the cost of the communications, while the ILP program of Section IV-F, so the comparison performed in Section V-D is more relevant than the one presented in [2].

The rest of this paper is organized as follows. Section II recalls the basics about Pareto dominance and how to compute the Pareto front with the  $\varepsilon$ -constraint method. Section III provides the required preliminaries including the application and architecture models and the interplay between the reliability, the power consumption, the temperature, and the execution time. Section IV provides the proposed scheduling heuristic ERPOT, along with its ILP counterpart. Section V presents the results of our simulations, performed both with syntactic benchmarks and with real-life benchmarks. Finally, Section VI surveys the related work and Section VII gives some concluding remarks.

## II. PARETO OPTIMIZATION

Before detailing our problem formulation, solutions, and algorithms, we give foundational background on Pareto optimization. When optimizing more than one criterion, there can be *several non-comparable* solutions, e.g., (42, 13) versus (9, 78) in the case of two criteria that must be minimized. The principle of Pareto optimization is to explore the design space by providing *as many solutions as possible*, to study the tradeoffs between these solutions. To compare solutions, we rely on the notion of *dominance* and Pareto optima, presented below in the case of two criteria that must be minimized (see Fig. 1(a)):

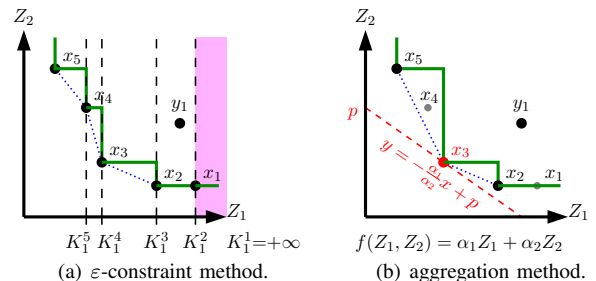


Fig. 1. Two transformation methods to compute the Pareto front (2D case).

- The point  $(x, y)$  *weakly dominates* the point  $(x', y')$  iff  $(x < x' \wedge y = y') \vee (x = x' \wedge y < y')$ . E.g.,  $x_2$  weakly dominates  $x_1$ .
- The point  $(x, y)$  *strongly dominates* the point  $(x', y')$  iff  $(x < x' \wedge y < y')$ . E.g.,  $x_3$  strongly dominates  $y_1$ .
- A point is a *weak Pareto optimum* iff there does not exist another point that strongly dominates it. E.g.,  $x_1, \dots, x_5$  are weak Pareto optima.
- A point is a *strong Pareto optimum* iff there does not exist another point that dominates it (weakly or strongly). E.g.,  $x_2, \dots, x_5$  are strong Pareto optima.
- The *Pareto front* is the set of all weak and strong Pareto optima.

Building the whole Pareto front and considering all constraints in a multi-criteria problem is a complicated task. To do this, several approaches exist [16], including the *aggregation* method that combines all the criteria in a single cost function, the *hierarchization* method that optimizes one criteria at a time, and the *transformation* method that transforms all the criteria except one into thresholds, and optimizes the remaining criterion under the constraints that the thresholds are satisfied (this last method is also called “budget optimization”). It is also possible to use population based methods, (e.g., genetic algorithms, particle swarm, ant colony, ...) or the Normal-Boundary Intersection method (NBI) [17].

Varying the cost function in the aggregation method or varying the order of the criteria in the hierarchization method can lead to computing several Pareto points, but not the entire Pareto front, a major theoretical drawback. The aggregation method is illustrated in Fig. 1(b) where the aggregation function is  $f(Z_1, Z_2) = \alpha_1 Z_1 + \alpha_2 Z_2$ . For two given values of  $\alpha_1$  and  $\alpha_2$ , the Pareto point that is found is the one that minimizes  $f$ : geometrically, it is the point from the Pareto front intersecting the line of slope  $-\alpha_1/\alpha_2$  and having the smallest value at origin (the  $p$  in Fig. 1(b)). The problem is that the concave portions of the Pareto front will be missed, e.g., the  $x_4$  point in Fig. 1(b); more generally, this is always the case if the aggregation function is convex (which is the case of  $f$ ). However, if the aggregation function is not convex, then there is no guarantee that the computed points are on the Pareto front. For instance, a non-convex aggregation function could return the point  $y_1$ .

Overall, the transformation method is an effective method to build the entire Pareto front when used in an iterative way. With two criteria, this is known as the  $\varepsilon$ -Constraint Method ( $\varepsilon$ CM) [12], depicted in Fig. 1(a). The criterion  $Z_1$  is transformed into a constraint. At iteration 1, the threshold for  $Z_1$  is set to  $K_1^1 = +\infty$ , yielding the Pareto optimum  $x_1$ . At iteration 2, the threshold  $K_1^2$  is set to the horizontal coordinate of  $x_1$ , therefore excluding the portion of the plane that is emphasized (in pink) and yielding the Pareto optimum  $x_2$ . This process repeats until all the points of the Pareto front have been found (if there is finite number of them), or until some pre-decided number of Pareto point have been found. Under the two conditions that (i) the number of Pareto optima is *finite* and that (ii) the minimization algorithm for  $Z_2$  computes the *optimal* result,  $\varepsilon$ CM computes the entire optimal Pareto front.

$\varepsilon$ CM has been later generalized to more than two criteria in [13], but at a very high computational cost:  $k^{m-1}\mathcal{O}(opt)$ , where  $k$  is the number of points in the Pareto front,  $m$  is the number of criteria, and  $\mathcal{O}(opt)$  is the complexity of the single criterion optimization algorithm. This computational complexity makes the generalized  $\varepsilon$ CM unfeasible for our problem (if only for the reason that the number of Pareto points is not bounded).

Instead, for each of the  $m-1$  criteria turned into a constraint, we simply divide the useful range of this criterion into  $p$  equally spaced intervals, and we invoke a single criterion optimization algorithm in each of the resulting  $p^{m-1}$  zones of the search space. We call this the Grid Method, depicted in Fig. 2(a) where the  $Z_1$  axis is divided into 4 intervals,

$[K_1^5, K_1^4)$  to  $[K_1^2, K_1^1)$ . The resulting complexity therefore becomes  $p^{m-1}\mathcal{O}(opt)$ . This is still exponential in  $m-1$  but the number of intervals  $p$  is much less than the number of Pareto points  $k$ . The number of intervals can be identical for each of the  $m-1$  criteria or not: each range can thus be divided into  $p_i$  intervals (not even necessarily equally spaced), resulting in an overall complexity of  $(\prod_{i=1}^{m-1} p_i)\mathcal{O}(opt)$ .

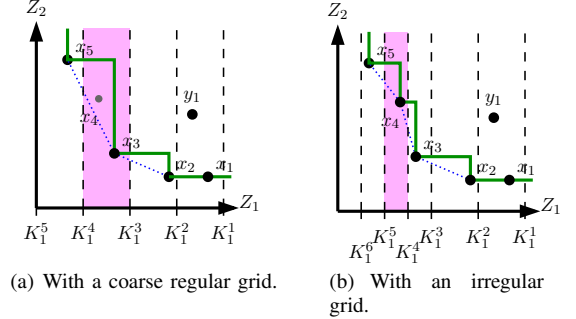


Fig. 2. The grid method to compute the Pareto front (2D case).

The choice of the intervals in each dimension has obviously an impact on the resulting Pareto front. For instance, the grid method illustrated in Fig 2(a) builds a Pareto front that does not include the point  $x_4$  because, in the interval  $[K_1^4, K_1^3)$  (emphasized in pink), the point that minimizes  $Z_2$  is  $x_3$ . With a different grid, the point  $x_4$  could be obtained, as shown in Fig. 2(b). On the one hand, using a finer grid will produce a Pareto front with more points, but this can become too costly. On the other hand, using an irregular grid could find more Pareto points but this seems very difficult to control a priori. It is simpler to generate the Pareto front with evenly spaced intervals in each dimension (the number of intervals depending on the time one is ready to spend to compute the Pareto front) and then, in order to improve locally the Pareto front around a particular Pareto optimum, to use either local search methods or to refine the intervals locally around this Pareto optimum. For instance, in Fig. 1, if  $x_3$  is identified as an interesting compromise, then the user can either use a local search algorithm around  $x_3$ , or he/she can divide the  $[K_1^4, K_1^3)$  interval into smaller intervals and invoke again the  $Z_2$  minimization function in these smaller intervals, which will be very likely to find the Pareto optimum  $x_4$ .

---

#### Algorithm 1 Grid method algorithm for 2 criteria.

---

**input:** The range  $[K_1^{min}, K_1^{max}]$  and the decrement  $\Delta$   
**output:** The list of Pareto points  $Res$   
1: **function** GRID( $K_1^{min}, K_1^{max}, \Delta$ )  
2:    $Res \leftarrow \emptyset$ ;  $K_1^i \leftarrow K_1^{max}$ ;  $i \leftarrow 1$   
3:   **while**  $K_1^i \geq K_1^{min}$  **do**  
4:      $Res \leftarrow Res \cup \text{OPT}(K_1^i)$   
5:      $K_1^i \leftarrow K_1^i - \Delta$   
6:   **end while**  
7:   **return** REMOVE\_NONDOMINATED\_POINTS( $Res$ )  
8: **end function**

---

To summarize, we use Algorithm 1 to implement the grid method, in the particular case of 2 criteria as in Fig. 2(b). The function  $\text{OPT}(K_1^i)$  returns the Pareto point that minimizes  $Z_2$

under the constraint  $Z_1 < K_1^i$ . The function REMOVE\_NON-DOMINATED\_POINTS( $Res$ ) removes the non dominated points from the list  $Res$  to produce the Pareto front.

As a final remark, note that in Fig. 1 and in Fig. 2, the Pareto front is depicted as a solid green line. It delimits the portion of the plane (above it and on its right) where all the points are dominated by a known Pareto optimum. This differs from the broken line that connects the Pareto optima (depicted in dotted blue), as is demonstrated by Fig. 1(b): the point  $x_4$  is above the dotted blue line, and yet we do not know whether or not it represents a feasible compromise between  $Z_1$  and  $Z_2$ , because no Pareto optimum has been found that dominates  $x_4$ .

### III. SYSTEM MODEL

#### A. Application and architecture models

An *application* is modeled as a directed acyclic graph (DAG)  $Alg = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V}$  is the set of nodes and  $\mathcal{E}$  is the set of edges. Each node represents a computing task, and each edge represents data-dependencies among two tasks. All tasks are assumed to be side-effect free (this assumption is required for active replication). If  $X \rightarrow Y$  is a data-dependency, then  $X$  is predecessor of  $Y$  and  $Y$  is successor of  $X$ .  $X$  is called the source of the data-dependency and  $Y$  is called its destination. We also define the sets  $pred(X) = \{Y | (Y, X) \in \mathcal{E}\}$  and  $succ(X) = \{Y | (X, Y) \in \mathcal{E}\}$ . Tasks with no predecessor are called *input* tasks, and those with no successors are called *output* tasks.

Fig. 3(a) shows an example of a DAG with two input tasks ( $I_1$  and  $I_2$ ), one output task ( $O_1$ ) and four regular tasks ( $A$ ,  $B$ ,  $C$  and  $D$ ).

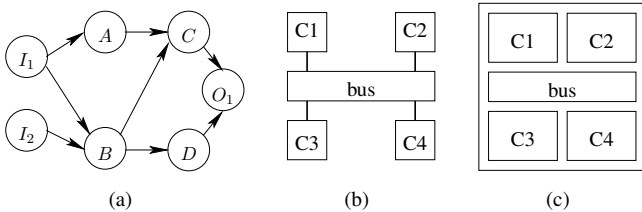


Fig. 3. (a) A sample application graph. (b) A sample architecture graph. (c) The corresponding coarse grain floorplan.

An *architecture* is a possibly heterogeneous multicore chip with one or more communication buses. It is modeled as a graph  $Arc = (\mathcal{C}, \mathcal{B}, \mathcal{L})$ , where  $\mathcal{C}$  is the set of cores,  $\mathcal{B}$  is the set of communication buses, and each  $e \in \mathcal{L}$  is a pair  $(c, b) \in \mathcal{C} \times \mathcal{B}$  specifying that the core  $c$  is connected to the bus  $b$ . We assume that there exists a path between any two cores  $c$  and  $c'$ . An example of a target architecture made of four cores and one bus is shown in Fig. 3(b).

We are also given a function  $\mathcal{E}xe_{nom}$  that returns the nominal (corresponding to the highest frequency) worst case execution times (WCETs) of all the tasks of  $Alg$  onto all the cores of  $Arc$ , as well as the worst case communication times (WCCTs) of all the data-dependencies of  $Alg$  onto all the communication buses of  $Arc$ . An intra-core communication takes no time to execute. For the sake of simplicity, all execution times are assumed to be integer numbers.

Computing the WCET of a given task on a processor has been the topic of much work. It involves finding the sequence of instructions in the program of the task that leads to the longest execution time. This is achieved by extracting the control flow graph (CFG) of the program, then by giving a duration (i.e., a number of clock cycles) to each basic block of the CFG. These durations are computed based on a model of the micro-architecture of the processor. This step includes some pessimism because of the hardware abstraction, be it in the cache replacement policy, the pipeline, the branch predictor, or the prefetch buffer. Based on this, the WCET is the length of the most weighted path in the annotated CFG. In general, the CFG contains backward edges, corresponding to the loops of the program. In this case, it is necessary to analyze the program in order to bound the number of iterations of each loop, which is classically done with abstract interpretation [18].

WCET analysis has been applied with success to real-life single-core processors actually used in embedded systems, with branch prediction [19] or with caches and pipelines [20]. These methods have later been adapted to multicores [21]–[23], taking into account the shared resources in the multicore (e.g., the shared memory or the bus).

Finally, the multicore is equipped with per-core DVFS. For each core, a set of (voltage, frequency) pairs  $\{(V_i, f_i)\}_{1 \leq i \leq \ell}$  is given. For the sake of simplicity, we assume that all the cores have the same set of (voltage, frequency) pairs. The actual execution time of a task  $\tau$  on a core  $c$  depends on the frequency  $f$  (in contrast, the buses are assumed to run at a fixed frequency denoted  $f_b$ ). To ease the computations, we transform the frequencies into scaling factors. E.g., if the set of available frequencies is  $\{900 MHz, 600 MHz, 300 MHz\}$ , then we use the scaling factors  $\{f_{max} = f_3 = 1, f_2 = \frac{2}{3}, f_{min} = f_1 = \frac{1}{3}\}$ . As a result, the WCET of task  $\tau$  at frequency  $f$  is given by:

$$\mathcal{E}xe(\tau, c, f) = \lceil \mathcal{E}xe_{nom}(\tau, c) / f \rceil \quad (1)$$

where the  $\lceil \cdot \rceil$  function guarantees that  $\mathcal{E}xe$  always returns an integer number.

#### B. Static mapping and scheduling

The specifications of the system consists of  $Alg$ ,  $Arc$ , and  $\mathcal{E}xe_{nom}$ . Implementing such a system involves two steps: First, we must find one or several cores of  $Arc$  to execute each task of  $Alg$ , and one or several communication buses of  $Arc$  for each data-dependency: this is the *mapping*. During this phase, we take into account (i) the reliability constraint by choosing how many cores must execute each task, (ii) the power consumption constraint by choosing at what frequency/voltage each component (core or bus) should execute each task and data-dependency, and (iii) the temperature constraint by inserting cooling times whenever necessary. Second, we must compute the starting time for each pair (task, proc) and each pair (data dep., bus): this is the *scheduling*. This paper solves these two steps *statically*, i.e., at compile time, based on a ready list scheduling heuristic. Finally, as said in the introduction, we schedule under constraints on the failure rate, the power consumption, and the temperature. We note respectively  $\Lambda_{obj}$ ,  $P_{obj}$ , and  $T_{obj}$  these constraints.



### C. Reliability

Both the cores and the buses are assumed to be *fail-silent*. Classically, we adopt the failure model of Shatz and Wang [24]: failures are *transient*, and the maximal duration of a failure is such that it affects only the current task executing onto the faulty core and not the subsequent tasks (same for the buses); this is known as the “hot” failure model.

Since the real-time systems we target are safety critical, the occurrence of failures is not acceptable and their reliability must be as close as possible to 1. One of the main causes of system failure are *transient failures* [25], which are commonly modeled by a Poisson distribution with a constant rate denoted  $\lambda$  [26]. Accordingly, the reliability of a single task or data-dependency  $\tau$  mapped onto a hardware component  $c$  (either a core or a bus) running at frequency  $f$  is:

$$R(\tau, c, f) = e^{-\lambda_c \cdot \mathcal{E}xe(\tau, c, f)} \quad (2)$$

where  $\lambda_c$  is the *failure rate per time unit* of the hardware component  $c$ , and  $\mathcal{E}xe(\tau, c, f)$  is the execution time of  $\tau$  on  $c$  at frequency  $f$ , computed with Eq. (1). When  $\tau$  is not replicated, we use Eq.(2). When  $\tau$  is actively replicated on a set  $\mathcal{K}$  of  $k$  hardware components numbered  $\{c_i\}_{1 \leq i \leq k}$ , each of them operating at frequency  $f_{c_i}$ , its reliability is:

$$R(\tau, \mathcal{K}) = 1 - \left( \prod_{i=1}^k \left( 1 - e^{-\lambda_{c_i} \cdot \mathcal{E}xe(\tau, c_i, f_{c_i})} \right) \right) \quad (3)$$

However, because of the operating frequency  $f$ ,  $\lambda$  is not constant anymore but is instead a function of the frequency [10]:

$$\lambda_f = \lambda_0 \cdot \rho_f \quad \text{with} \quad \rho_f = 10^{\frac{b(1-f)}{1-f_{min}}} \quad (4)$$

where  $\lambda_0$  is the nominal failure rate per time unit,  $\rho_f$  is the frequency-dependent factor,  $b$  is a strictly positive constant that accounts for the susceptibility of hardware to transient faults due to frequency scaling,  $f$  is the operational frequency level, and  $f_{min}$  is the lowest frequency of the system. Recall that the frequency value  $f$  is normalized in the range  $[0, 1]$  with  $f_{max} = 1$ . This is consistent with Eq (1).

Many articles have studied the impact of the temperature on the rate of transient faults [27]–[29]. In addition, there are several mechanisms that lead to permanent failures, most notably electro-migration, negative bias temperature instability, stress migration, time-dependent dielectric breakdown, and thermal cycling [3], [30]. All of these phenomena can be characterized by a failure rate as an exponential function of the temperature. We take into account the effect of the temperature on the failure rate per time unit with the Arrhenius equation [3]:

$$\lambda_T = \lambda_0 \cdot \rho_T \quad \text{with} \quad \rho_T = e^{-\frac{E_a}{K} \cdot \left( \frac{1}{T(t)} - \frac{1}{T_0} \right)} \quad (5)$$

where again  $\lambda_0$  is the nominal failure rate per time unit,  $\rho_T$  is the temperature-related factor,  $E_a$  is the activation energy,  $K$  is the Boltzmann’s constant,  $T(t)$  is the temperature of the system at time  $t$  in Kelvin, and  $T_0$  is the initial temperature. Of course, we will also have to take into account the effect of each core’s temperature on the other cores (see Section III-E).

Finally, we combine Eqs (4) and (5) to provide a global equation of the failure rate per time unit as a function of the frequency and the temperature. Since the frequency factor  $\rho_f$  and the temperature factor  $\rho_T$  are both dimension-less, the dimension of  $\lambda_{sys}$  is the same as  $\lambda_0$ , hence  $\lambda_{sys}$  is also a failure rate per time unit:

$$\lambda_{sys} = \lambda_0 \cdot \rho_f \cdot \rho_T = \lambda_0 \cdot 10^{\frac{b(1-f)}{1-f_{min}}} \cdot e^{-\frac{E_a}{K} \cdot \left( \frac{1}{T(t)} - \frac{1}{T_0} \right)} \quad (6)$$

When computing the reliability of a given task or data-dependency  $\tau$  on a single hardware component  $c_i$  (resp. a set  $\mathcal{K} = \{c_i\}_{1 \leq i \leq k}$ ), we therefore use Eq. (2) (resp. Eq. (3)) by replacing  $\lambda_{c_i}$  by  $\lambda_{sys}(c_i)$ :

$$R(\tau, c_i, f_{c_i}, t) = e^{-\lambda_{sys}(c_i) \cdot \mathcal{E}xe(\tau, c_i, f_{c_i})} \quad (7)$$

$$R(\tau, \mathcal{K}, t) = 1 - \left( \prod_{i=1}^k \left( 1 - e^{-\lambda_{sys}(c_i) \cdot \mathcal{E}xe(\tau, c_i, f_{c_i})} \right) \right) \quad (8)$$

where  $t$  is shown to make explicit the dependency of the temperature of  $c_i$  on the time in  $\lambda_{sys}(c_i)$ . In the entire paper, we take the temperature at the *task granularity*, i.e., we assume that  $T(t)$  remains constant for the entire duration of  $\tau$ . We will prove at the end of this section that doing so is safe for the  $\Lambda_{obj}$  constraint.

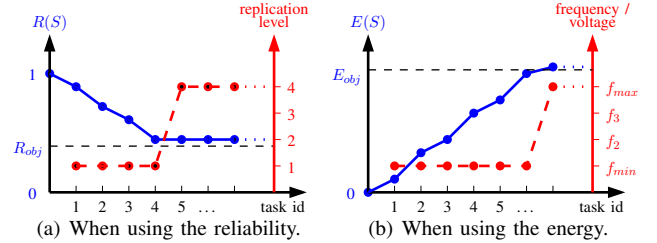


Fig. 4. Funnel effect.

It has been demonstrated in [2], [5] that using the reliability as a constraint in the  $\varepsilon$ -constraint method *does not work*. Intuitively, this is because the reliability is not an invariant measure of the number of scheduled tasks. Indeed, computing the reliability of a schedule involves, at each mapping decision, a multiplication by a factor that is strictly less than 1: see Eq.(3). This is illustrated in Fig. 4(a), where the horizontal axis counts the task numbers in their mapping order (recall that we use a ready list scheduling algorithm). As long as the reliability is above the threshold  $R_{obj}$ , the tasks are not replicated, because this is what minimizes the schedule length; thus the replication level of tasks 1 to 4 is 1 (red dashed line). This results in a multiplicative factor significantly below 1, which causes the system’s reliability to drop (blue solid line). Once task 4 has been scheduled, the reliability is very close to  $R_{obj}$ ; this causes the replication level to skyrocket up to a value sufficient for the multiplying factor to be close enough to 1, so that the system’s reliability remains above  $R_{obj}$ . We call this the “funnel effect” [2].

For this reason, instead of the reliability, we use the *Global System Failure Rate* (GSFR) [5]. Intuitively, the GSFR of a possibly partial schedule is the failure rate of the system operating under this schedule *as if* it was a single task mapped

on a single core. As a consequence, we schedule under a constraint  $\Lambda_{obj}$  on the GSFR instead of a constraint  $R_{obj}$  on the reliability. For a single task  $\tau$ , the GSFR is denoted  $\Lambda(\tau)$  and is computed as:

$$\Lambda(\tau, c, f, t) = \frac{-\log(R(\tau, c, f, t))}{\mathcal{E}xe(\tau, c, f)} \quad (9)$$

And for a schedule  $S$ , the GSFR  $\Lambda(S)$  is computed as:

$$\Lambda(S) = \frac{-\log(R(S))}{U(S)} \quad \text{with} \quad U(S) = \sum_{(\tau, c, f) \in S} \mathcal{E}xe(\tau, c, f) \quad (10)$$

where  $R(S)$  is the reliability of the schedule  $S$  and  $U(S)$  is the overall sum of the execution times of the cores in  $S$ . The notation  $(\tau, c, f) \in S$  means that, in the schedule  $S$ , task  $\tau$  is executed on core  $c$  at frequency  $f$ . Eq.(10) is equivalent to  $R(S) = e^{-\Lambda(S) \cdot U(S)}$ , which is the same as Eq. (2) but for a schedule  $S$  instead of a single task  $\tau$ .

One key aspect of Eq. (10) is that it uses  $U(S)$  and not the schedule length. There are two reasons behind this choice: first it makes the computation of the GSFR *compositional* with respect to the structure of the schedule, and second it is consistent with the ‘‘hot’’ failure model [5].

The consequence of this shift from the reliability to the GSFR is that, from now on, our state space will be the 4D space (execution time, GSFR, power, temperature).

We are now ready to prove that assuming the temperature on each core  $c_j$  and on the bus  $b$  to remain constant during the duration of each task/data-dependency  $\tau$  is safe w.r.t. the  $\Lambda_{obj}$  constraint.

**Proposition 1:** Let  $\tau$  be a task or a data-dependency scheduled on a hardware component  $c$  at frequency  $f$ , starting at time  $t_0$  and finishing at time  $t_f = t_0 + \mathcal{E}xe(\tau, c, f)$ . The reliability of  $\tau$  on  $c$  is computed with Eq. (7) and the GSFR with Eq. (9). (i) If the temperature increases over the interval  $[t_0, t_f]$ , then fixing  $T(t) = T(t_f)$  is safe regarding the  $\Lambda_{obj}$  constraint. (ii) If the temperature decreases over the interval  $[t_0, t_f]$ , then fixing  $T(t) = T(t_0)$  is safe regarding the  $\Lambda_{obj}$  constraint.

**Proof:** (i) In the heating mode, the temperature increases during the execution of  $\tau$ , and when it does,  $\lambda_{sys}(c)$  increases too. Since  $R$  is decreasing in function of  $\lambda_{sys}$ , we have:

$$\forall t \in [t_0, t_f], R(\tau, c, f, t) \geq R(\tau, c, f, t_f)$$

Since  $R(\tau, c, f, t) \geq R(\tau, c, f, t_f) \iff \Lambda(\tau, c, f, t) \leq \Lambda(\tau, c, f, t_f)$ , we therefore have:

$$\Lambda(\tau, c, f, t_f) \leq \Lambda_{obj} \implies \forall t \in [t_0, t_f], \Lambda(\tau, c, f, t) \leq \Lambda_{obj}$$

which proves that assuming that  $T(t)$  remains constant and equal to  $T(t_f)$  is safe regarding the  $\Lambda_{obj}$  constraint.

(ii) In the cooling mode, the proof is identical since  $T(t)$  decreases so  $\lambda_{sys}(t)$  decreases, hence assuming that  $T(t)$  remains constant and equal to  $T(t_0)$  is safe regarding the  $\Lambda_{obj}$  constraint.  $\square$

#### D. Power consumption

The power consumption of a single task (or data-dependency) running on a hardware component is composed of two aspects [10], [31]: (i) the leakage power and (ii) the dynamic power. The former depends on the leakage current, which itself mostly depends on the chip temperature, while the latter depends on the chosen pair (voltage  $V$ , frequency  $f$ ). The overall power consumption  $P_{sys}$  is equal to  $P_{leak} + P_{dyn}$ , computed by Eq. (11):

$$\begin{cases} P_{sys}(t) = \alpha \cdot T(t) + \beta_h + \gamma \cdot C_{ef} \cdot V^2 \cdot f & \text{if heating} \\ P_{sys}(t) = \alpha \cdot T(t) + \beta_c + \gamma \cdot C_{ef} \cdot V^2 \cdot f & \text{if cooling} \end{cases} \quad (11)$$

Regarding the leakage power,  $\alpha$ ,  $\beta_h$ , and  $\beta_c$  are architecture-dependent coefficients and are determined based on the characteristics of the platform;  $\beta_h$  is used in the heating mode and  $\beta_c$  in the cooling mode [8]. Finally,  $T(t)$  is the chip temperature at time  $t$ , in Kelvin. Regarding the dynamic power,  $V$  is the supply voltage,  $f$  is the frequency,  $C_{ef}$  is the switching capacitance (a constant that depends on the chip technology), and  $\gamma$  is the activity ratio, which varies from 0 (no activity) to 1 (all gates are active at each cycle). In theory, there should be a different  $\gamma$  for each task, and our scheduling algorithm can handle it. In practice, for the sake of simplicity we take an average  $\gamma$  value, identical for all the tasks.

Recall that we take the temperature at the *task granularity*, i.e., we assume that  $T(t)$  remains constant for the entire duration of  $\tau$ . The following property states that doing this is safe regarding the  $P_{obj}$  constraint.

**Proposition 2:** Let  $\tau$  be a task or a data-dependency scheduled on a hardware component  $c$  at frequency  $f$ , starting at time  $t_0$  and finishing at time  $t_f = t_0 + \mathcal{E}xe(\tau, c, f)$ . The power consumption of  $c$  during the execution of  $\tau$  is computed with Eq. (11). (i) If the temperature increases over the interval  $[t_0, t_f]$ , then fixing  $T(t) = T(t_f)$  is safe regarding the  $P_{obj}$  constraint. (ii) If the temperature decreases over the interval  $[t_0, t_f]$ , then fixing  $T(t) = T(t_0)$  is safe regarding the  $P_{obj}$  constraint.

**Proof:** (i) In the heating mode, the temperature increases during the execution of  $\tau$ , and when it does,  $P_{sys}(t)$  increases too. It follows that assuming  $T(t)$  to remain constant over the interval  $[t_0, t_f]$  and equal to  $T(t_f)$  yields  $\forall t, P_{sys}(t) \leq P_{sys}(t_f)$ . Therefore, we have:

$$P_{sys}(t_f) \leq P_{obj} \implies \forall t \in [t_0, t_f], P_{sys}(t) \leq P_{obj}$$

which proves that assuming that  $T(t)$  remains constant and equal to  $T(t_f)$  is safe regarding the  $P_{obj}$  constraint.

(ii) In the cooling mode, the proof is identical since  $T(t)$  decreases so  $P_{sys}(t)$  decreases, hence assuming that  $T(t)$  remains constant and equal to  $T(t_0)$  is safe regarding the  $P_{obj}$  constraint.  $\square$

From Eq. (11), we can then compute the energy consumed by the system when executing a schedule (possibly partial). However, the same funnel effect as with the reliability occurs if one uses the energy as a constraint in the  $\varepsilon$ -constraint method [2]. The reason again is that the energy is not an invariant measure of the number of scheduled tasks. Indeed,

computing the energy consumed by a schedule involves, at each mapping decision, an addition of a term that is strictly positive. This is illustrated in Fig. 4(b): the horizontal axis counts the task numbers in their mapping order; the blue solid line depicts the cumulative energy consumed by the system; up to task 6, the energy is below the energy constraint  $E_{obj}$  so everything is fine; however, there is no possibility to schedule task 7 without violating the energy constraint. For this reason, in our multi-criteria scheduling heuristic we use the power consumption, with a constraint  $P_{obj}$ , which is an invariant measure of the number of scheduled tasks.

### E. Temperature

The instantaneous temperature of a computing system depends on the power consumption and on the current temperature (and its variations in time). For a given hardware component  $c$  (core or bus), it is computed based on the following differential equation [32]:

$$C \cdot \left( \frac{dT_c(t)}{dt} \right) + G(T_c(t) - T_{amb}) = P(t) \quad (12)$$

where  $C$  and  $G$  are the architecture-based constants for the heat conductivity,  $T_c$ ,  $t$ ,  $T_{amb}$ , and  $P$  are respectively the temperature of  $c$ , the time, the ambient temperature (assumed to be less than  $T_{obj}$ <sup>1</sup>), and the instantaneous power consumption of the system. The power consumption is the sum of the static and dynamic power, as given by Eq. (11).

For each component  $c$ , we wish to take into account the effect of the temperature of its neighbors, according to the chip floorplan. To achieve this, we add a heat transfer term to Eq. (12):

$$C \cdot \left( \frac{dT_c(t)}{dt} \right) + G(T_c(t) - T_{amb}) + 2D_{heat} = P(t) \quad (13)$$

and we use the coarse grain floorplan of Fig. 3(c) (similar to the spatial thermal model and floorplan of [33]) to model this two-dimension heat transfer as:

$$2D_{heat} = \sum_{c' \in nbr(c)} \kappa(c, c') \cdot (T_c(t) - T_{c'}(t)) \quad (14)$$

where  $nbr(c)$  is the set of all neighbors of  $c$ ,  $T_c$  is the temperature of  $c$ , and  $\kappa(c, c')$  is the thermal conductivity between  $c$  and  $c'$ , which depends on their distance and on the chip geometry characteristics (as given in the floorplan).

Combining Eqs. (13) and (14) yields the following differential equation for  $T_c(t)$ :

$$\begin{aligned} C \cdot \left( \frac{dT_c(t)}{dt} \right) &= -G \cdot (T_c(t) - T_{amb}) \\ &- \sum_{c' \in nbr(c)} \kappa(c, c') \cdot (T_c(t) - T_{c'}(t)) \\ &+ C_{ef} \cdot V^2 \cdot f + \alpha \cdot T_c(t) + \beta_h \end{aligned} \quad (15)$$

<sup>1</sup>If  $T_{amb} > T_{obj}$ , then putting the component in the idle mode does not allow it to cool down.

We then proceed as in [8] to re-write Eq. (15) as:

$$\begin{aligned} \frac{dT_c(t)}{dt} &= -A \cdot T_c(t) + B, \quad \text{with} \\ A &= \frac{G - \alpha + \sum_{c' \in nbr(c)} \kappa(c, c')}{C} \\ B &= \frac{G \cdot T_{amb} + \sum_{c' \in nbr(c)} \kappa(c, c') \cdot T_{c'}(t) + C_{ef} \cdot V^2 \cdot f + \beta_h}{C} \end{aligned} \quad (16)$$

When computing the evolution of the temperature of  $c$  during the execution of  $\tau$ , we assume that the temperatures of the neighbors remain constant for the entire duration of  $\tau$ , and equal to their respective temperature at the end of  $\tau$ . By virtue of the same reasoning as the one made in Section III-D, this is safe regarding the  $T_{obj}$  constraint. It follows that the closed form solution to the differential equation (16) is:

$$T_c(t) = T_{\infty}^{heat} + (T_0 - T_{\infty}^{heat}) \cdot e^{-A(t-t_0)} \quad (17)$$

where  $T_{\infty}^{heat} = B/A$  is the heating steady state temperature and  $T_0 = T(t_0)$  is the temperature of the system at  $t_0$ . We note  $T_c^{heat}(t_0, t)$  the temperature computed with Eq. (17).

When  $c$  is idle, the computation of the temperature is identical except that the term  $C_{ef} \cdot V^2 \cdot f$  in Eq. (11) disappears and  $\beta_h$  is replaced by  $\beta_c$ , yielding the following closed form:

$$\begin{aligned} B' &= \frac{G \cdot T_{amb} + \sum_{c' \in nbr(c)} T_{c'}(t) \cdot \kappa(c, c') + \beta_c}{C} \\ T_c(t) &= T_{\infty}^{cool} + (T_0 - T_{\infty}^{cool}) \cdot e^{-A(t-t_0)} \end{aligned} \quad (18)$$

where  $T_{\infty}^{cool} = B'/A$  is the cooling steady state temperature and  $T_0 = T(t_0)$  is the temperature of the system at  $t_0$ , i.e., at the start of the cooling time. We note  $T_c^{cool}(t_0, t)$  the temperature computed with Eq. (18).

## IV. ERPOT: THE PROPOSED QUAD-CRITERIA OPTIMIZATION SCHEDULING HEURISTIC METHOD

The optimal mapping of a DAG of tasks on a multicore is a known NP-complete problem [15]. We therefore propose a heuristic algorithm, more precisely a ready list scheduling, for which we formally prove that each computed schedule satisfies the GSFR, power consumption, and temperature constraints (Section IV-B). In addition, in order to assess the performances of our heuristic, we implement an optimal version on top of an ILP solver (Section IV-F).

### A. General principles of ERPOT

We are given:

- (i) a DAG of tasks  $Alg = (\mathcal{V}, \mathcal{E})$ ,
- (ii) a multicore architecture description  $Arc = (\mathcal{C}, \mathcal{B}, \mathcal{L})$  along with the nominal failure rate per time unit  $\lambda_0$  of each hardware component,



- (iii) a function  $\mathcal{E}xe_{nom}$  of the nominal WCETs / WCCTs of of all the tasks / data-dependencies of  $\mathcal{Alg}$  onto all the cores / buses of  $\mathcal{Arc}$ ,
- (iv) a set of frequencies for the cores  $\mathcal{F} = \{f_j\}_{1 \leq j \leq \ell}$  and a fixed frequency  $f_b$  for the buses, all taken as scaling factors,
- (v) three constraints  $\Lambda_{obj}$ ,  $P_{obj}$ , and  $T_{obj}$  respectively on the GSFR, the power consumption, and the temperature,
- (vi) and the initial temperature of the chip  $T_{init}$ .

The goal is to compute, if it exists, a schedule of  $\mathcal{Alg}$  onto  $\mathcal{Arc}$  such that the three constraints are met and the execution time is minimal<sup>2</sup>. If no solution is found, it means that the available hardware resources are not sufficient to meet the desired constraints  $\Lambda_{obj}$ ,  $P_{obj}$ , and  $T_{obj}$ . This issue is discussed in Section IV-B.

In order to keep the GSFR below  $\Lambda_{obj}$ , we use the *active replication* of tasks. We compute the reliability of a partial schedule by building the corresponding Reliability Block Diagram (RBD) [26]. An RBD is a DAG that starts with a source node  $S$  and ends with a destination node  $D$ . Between  $S$  and  $D$ , each of its nodes corresponds to one task (or data-dependency) scheduled on a core (or bus). By definition, an RBD is *operational* iff there exists at least one operational path from  $S$  to  $D$ . A path is operational iff all the blocks in this path are operational. The probability that a block is operational is its reliability, computed with Eq. (2). By construction, the probability that an RBD is operational is therefore equal to the reliability of the static schedule it represents.

Computing the reliability in this way assumes that the occurrences of the failures are *statistically independent events*. Without this hypothesis, the fact that some blocks belong to several paths from  $S$  to  $D$  makes the computation of the reliability very complex. Concerning hardware faults, this hypothesis is reasonable, but this would not be the case for software faults [34].

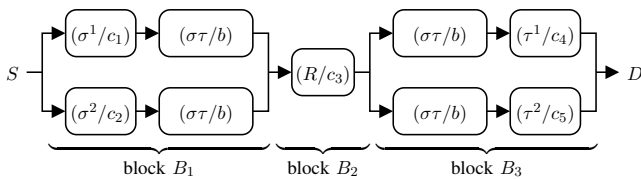


Fig. 5. Reliability Block Diagram for a simple schedule with replication.

In general, the structure of the RBD is unspecified, which makes the reliability computation NP-complete [35]. Following [5], the solution we use to prevent this is to insert *routing* tasks (the execution time of which is 0) from each set of replicas of a predecessor task to the set of replicas of its successor task. As a result, the RBD is *serial-parallel*, which makes the reliability computation linear. For any task  $\tau$  of  $\mathcal{Alg}$ , all its replicas appear in parallel in the same block of the RBD, whose reliability is therefore computed by Eq. (3), and the RBD is composed of all these blocks in sequence (hence the serial-parallel structure). Consider for instance a simple

DAG with two tasks  $X \rightarrow Y$  to be scheduled onto a six-core chip with a single bus. If  $X$  is replicated twice on cores  $c_1$  and  $c_2$  (its two replicas being denoted  $X^1$  and  $X^2$ ), and  $Y$  is replicated twice on cores  $c_4$  and  $c_5$  (its two replicas being denoted  $Y^1$  and  $Y^2$ ), then the RBD for this schedule will have the form shown in Fig. 5.

In practice, our scheduling heuristics will optimize the placement of the routing tasks so as to minimize the total execution time, for instance by mapping  $R$  to  $c_1$  or  $c_2$ .

Owing to the serial-parallel structure of the RBD, computing the reliability of a schedule (be it partial or final) is compositional. It follows that, to guarantee that the entire schedule satisfies the  $\Lambda_{obj}$  constraint, it suffices to guarantee that each block of the RBD satisfies this constraint.

In order to keep the power consumption below  $P_{obj}$ , we use two techniques: (i) on the one hand DVFS, which is available on many modern multicores such as the Intel i7-2600 quad-core or the Samsung Exynos 5422 octa-core; this allows us to lower the  $P_{dyn}$  term of Eq. (11); and (ii) on the other hand we try to keep the temperature below  $T_{obj}$ , which allows us to lower the  $P_{leak}$  term of Eq. (11). Computing the dynamic power consumption requires computing the energy consumed by the schedule (be it partial or local), and then to divide it by the schedule length. The compositionality issue raised by the GSFR computation also arises here. As demonstrated in [2], this issue can be solved by *over-estimating* the energy consumption each time that the partial schedule has a “hole” at the end, that is, each time one of the cores is idle while the other cores are busy executing their last task. Over-estimation is achieved by computing the energy consumed by such a schedule as if the “hole” was “filled” with a virtual task running at the maximal frequency.

In order to keep the temperature below  $T_{obj}$ , we insert *cooling times* to allow the cores to cool down [8], [36], [37] (the buses are always much less loaded than the cores, so they never need to cool down). We follow the same principle as the JUST strategy proposed in [8] for single-core processors, with two differences: first, the target architecture is a multicore, and second, our objective is to minimize the schedule length under a maximal temperature constraint. The rationale of the JUST strategy is to insert cooling times as late as possible and only when needed, i.e., just in time. Thus, each time we want to schedule a task  $\tau$  on a core  $c$ , we evaluate the temperature of each core in the multicore at the end of this task, taking into account the planned voltage and frequency of  $\tau$  and the influence of the temperature of the neighbors of  $c$ . If it exceeds  $T_{obj}$ , then we postpone the starting time of  $\tau$  by inserting a *cooling time* in order to cool down the core  $c$ . The length of the cooling time is the *smallest length* such that the temperature at the end of  $\tau$  does not exceed  $T_{obj}$ .

Recall that a high temperature has a negative effect on the reliability (as shown in Eq. (6)) as well as on the leakage power consumption (see Eq. (11)). This makes it all the more important to limit the maximal temperature.

### B. Quad-criteria scheduling heuristic algorithm

ERPOT is a ready list scheduling algorithm implemented in MATLAB (1,300 lines of code). It works with two lists, the list

<sup>2</sup>The execution time of the schedule is also called  $C_{max}$  or schedule length in the scheduling community.

$\mathcal{Ready}^{(n)}$  of ready tasks and the list  $\mathcal{Sched}^{(n)}$  of scheduled tasks, where  $(n)$  denotes the current step of the list scheduling. At each step  $(n)$ , we have  $\mathcal{Ready}^{(n)} \cap \mathcal{Sched}^{(n)} = \emptyset$ .

In a preliminary phase, we traverse the  $\mathcal{Alg}$  graph breadth-first, from the output tasks to the input tasks, in order to compute, for each task  $\tau$ , the *Longest Execution Path* from  $\tau$  to the end of the graph, noted  $LEP(\tau)$ . This notion is similar to the “bottom-level” presented in [38]. Intuitively,  $LEP(\tau)$  accounts for all the “future” tasks of  $\tau$ . For each task  $\tau$ , it is computed as follows:

- If  $\text{succ}(\tau) = \emptyset$ , then we compute its  $LEP$  as  $LEP(\tau) = (\sum_{c \in \mathcal{C}} \mathcal{E}xe_{nom}(\tau, c))/|\mathcal{C}|$ . The nominal execution time of  $\tau$  is averaged over all the cores (the set  $\mathcal{C}$ ) since we do not know in advance onto which core  $\tau$  will be actually scheduled.
- If  $\text{succ}(\tau) = \{\tau'\}$ , then  $LEP(\tau) = LEP(\tau') + (\sum_{c \in \mathcal{C}} \mathcal{E}xe_{nom}(\tau, c))/|\mathcal{C}|$ . Since  $\tau$  has only one successor, its nominal execution time is added to the  $LEP$  of its only successor (again, averaged over all cores).
- If  $\text{succ}(\tau) = \{\tau_i\}_{1 \leq i \leq k}$  with  $k \geq 2$ , then  $LEP(\tau) = \max_{1 \leq i \leq k} LEP(\tau_i) + (\sum_{c \in \mathcal{C}} \mathcal{E}xe_{nom}(\tau, c))/|\mathcal{C}|$ . Since  $\tau$  has more than one successor, its averaged nominal execution time is added to the max of the  $LEPs$  of all its successors (again, averaged over all cores).

Still in the preliminary phase, we build the set  $2^{\mathcal{C}}$  of all subsets of  $\mathcal{C}$ , and for each such subset  $\{c_i\}_{1 \leq i \leq k} \in 2^{\mathcal{C}}$ , we build all the possible sets of pairs  $\{(c_i, f_j)\}_{1 \leq i \leq k, 1 \leq j \leq \ell}$ , where  $\ell$  is the number of available frequencies. We denote by  $\mathcal{Q}$  the set of all such sets of pairs (core, frequency).

In the main phase of ERPOT, we first assign to  $\mathcal{Ready}^{(0)}$  the set of input tasks of  $\mathcal{V}$ , and to  $\mathcal{Sched}^{(0)}$  the empty set. Then, at each step  $(n)$ , we select the *most urgent* task to be scheduled among all the ready tasks, that is, the task  $\tau_{urg}$  for which  $LEP(\tau)$  is the largest:  $\tau_{urg} = \text{argmax}_{\tau \in \mathcal{Ready}^{(n)}} LEP(\tau)$ .

The next step involves selecting the *best subset of cores* and their associated frequencies to execute  $\tau_{urg}$ . Each  $Q_i \in \mathcal{Q}$  is a potential scheduling choice for  $\tau_{urg}$ , which we need to evaluate according to our three constraints and our minimization criterion. We denote by  $Q_{best}$  the best scheduling choice, by  $L^{(n)}$  the schedule length at step  $(n)$ , thus *before* executing  $\tau$  on  $Q_{best}$ , and by  $L^{(n+1)}(\tau, Q_{best})$  the schedule length *after* executing  $\tau$  on  $Q_{best}$ , which we shorten into  $L^{(n+1)}$  to avoid heavy notations. Similarly, we denote by  $\Lambda^{(n)}$  the GSFR,  $E^{(n)}$  the energy, and  $T^{(n)}$  the temperature at step  $(n)$ , again shortened. We further note  $\Lambda(\tau, Q_{best})$  the GSFR of the parallel block corresponding to executing  $\tau$  onto each core of  $Q_{best}$ . Recall that we have explained in Section IV-A that the GSFR of a schedule is computed block by block, as a result of the serial-parallel structure of its RBD. With these notations,  $Q_{best}$  is given by the following equation:

$$Q_{best} = \underset{Q_i \subseteq \mathcal{Q}}{\text{argmin}} \left\{ \begin{aligned} & \Lambda(\tau, Q_i) \leq \Lambda_{obj} \\ & \wedge (E^{(n+1)} - E^{(n)}) \leq P_{obj} \cdot (L^{(n+1)} - L^{(n)}) \\ & \wedge T^{(n+1)} \leq T_{obj} \end{aligned} \right\} \quad (19)$$

Eq. (19) might return an empty set  $Q_{best}$ . This can occur for three reasons:

- 1) Either there is no subset of cores  $Q_i$  that satisfies the GSFR criterion  $\Lambda(\tau, Q_i) \leq \Lambda_{obj}$ . In other words, the number of available cores is not sufficient to reach the required GSFR level. The heuristic fails and returns a “no solution” result. Recall that we want to find solutions in the 4D space (execution time, GSFR, power, temperature). So “no solution” only means that there will be no Pareto point at the coordinates  $(\Lambda_{obj}, P_{obj}, T_{obj})$  in the 4D space.
- 2) Either there is no subset of cores  $Q_i$  that satisfies the power consumption criterion  $(E^{(n+1)} - E^{(n)}) \leq P_{obj} \cdot (L^{(n+1)} - L^{(n)})$ . In other words, the available frequencies are not sufficient to reach the required power consumption level. Like in case 1 above, the heuristic fails and returns a “no solution” result.
- 3) Or there is no subset of cores  $Q_i$  that satisfies the temperature criterion  $T^{(n+1)} \leq T_{obj}$ . In this case, let  $Q'_i = \{c_j \in Q_i \mid T^{(n+1)}(c_j) > T_{obj}\}$  and let  $t_j$  be earliest time at which  $\tau$  can start on core  $c_j$ . We add to each core  $c_j \in Q'_i$  a cooling time of length  $s_j$  that starts at  $t_j$ , such that  $s_j$  is the smallest integer satisfying the inequality:

$$T_{c_j}^{cool}(t_j, s_j) + T_{c_j}^{heat}(t_j + s_j, \mathcal{E}xe(\tau, c_j, f_j)) \leq T_{obj}$$

### C. Soundness of our scheduling heuristic

We prove in this section four key propositions on the produced schedules, which guarantee that the schedules generated by ERPOT satisfy the  $\Lambda_{obj}$ ,  $P_{obj}$ , and  $T_{obj}$  constraints.

**Proposition 3:** Let  $S$  be a schedule of  $\mathcal{Alg}$  onto  $\mathcal{Arc}$ . If each task of  $\mathcal{Alg}$  has been scheduled on the subset of cores  $Q_{best}$  defined by Eq. (19), thus satisfying the GSFR constraint  $\Lambda_{obj}$ , then the total schedule  $S$  will also meet the  $\Lambda_{obj}$  constraint.

**Proof (see [5]):** Each task  $\tau_i$  of  $\mathcal{Alg}$  is scheduled onto a subset  $Q_{best}^i$  that was selected by Eq. (19). Hence, for all  $\tau_i$  in  $\mathcal{Alg}$ , we have  $\Lambda(\tau_i, Q_{best}^i) \leq \Lambda_{obj}$ . Owing to the serial-parallel structure of the RBD corresponding to the schedule  $S$  and to the fact that the GSFR is computed compositionally from the RBD, it follows that  $\Lambda(S) \leq \Lambda_{obj}$ .  $\square$

**Proposition 4:** Let  $S$  be a schedule of  $\mathcal{Alg}$  onto  $\mathcal{Arc}$ . If each task of  $\mathcal{Alg}$  has been scheduled on the subset of cores  $Q_{best}$  defined by Eq. (19), thus satisfying the power consumption constraint  $P_{obj}$ , then the total schedule  $S$  will also meet the  $P_{obj}$  constraint.

**Proof (see [2]):** The proof follows from Eq. (19) and from the compositionality of the power consumption (as opposed to the energy). Notice that the constraint on the power consumption in Eq. (19) is actually expressed as a constraint between the *energy increase*  $(E^{(n+1)} - E^{(n)})$  and the *schedule length increase*  $(L^{(n+1)} - L^{(n)})$ . The reason is the following: suppose that, at step  $(n)$ , the most urgent task is  $\tau_i$  with  $Q_{best}^i = \{(c_i, f_i)\}$ ; suppose also that, in the partial schedule before mapping  $\tau_i$ , the finish time  $L_{c_i}$  on core  $c_i$  is such that  $L_{c_i} + \mathcal{E}xe(\tau_i, c_i, f_i) \leq L^{(n)}$ ; in other words, scheduling  $\tau_i$  on  $c_i$  at frequency  $f_i$  does *not* increase the current schedule length because there is a “hole” at the end of the schedule of core  $c_i$ . Hence  $L^{(n)} = L^{(n+1)}$ . In contrast, the energy does increase when  $\tau_i$  is scheduled on  $c_i$  at frequency  $f_i$ , so  $E^{(n+1)} > E^{(n)}$ .

To overcome this issue, we have proposed in [2] a solution where we “fill” each “hole” at the end of the schedule with a virtual task executing at the maximal frequency  $f_{max}$ . It follows the energy consumed by the partial schedule at each step ( $n$ ) is *over-estimated*.

With this over-estimation, we prove the desired property by induction on ( $n$ ). At step (1), the property is verified because the first task is scheduled according to Eq. (19):

$$(E^{(1)} - E^{(0)}) \leq P_{obj} \cdot (L^{(1)} - L^{(0)}) \iff P^{(1)} = \frac{E^{(1)}}{L^{(1)}} \leq P_{obj}$$

Then, our induction hypothesis is:

$$P^{(n)} = \frac{E^{(n)}}{L^{(n)}} \leq P_{obj} \iff E^{(n)} \leq P_{obj} L^{(n)} \quad (20)$$

As a result of Eq. (19), we have:

$$\begin{aligned} (E^{(n+1)} - E^{(n)}) &\leq P_{obj} \cdot (L^{(n+1)} - L^{(n)}) \\ \iff E^{(n+1)} &\leq E^{(n)} + P_{obj} L^{(n+1)} - P_{obj} L^{(n)} \end{aligned}$$

Owing to the induction hypothesis (20), this implies:

$$\begin{aligned} E^{(n+1)} &\leq P_{obj} L^{(n)} + P_{obj} L^{(n+1)} - P_{obj} L^{(n)} \\ \iff E^{(n+1)} &\leq P_{obj} L^{(n+1)} \\ \iff P^{(n+1)} &= \frac{E^{(n+1)}}{L^{(n+1)}} \leq P_{obj} \end{aligned}$$

which concludes the proof by induction.  $\square$

**Proposition 5:** Let  $S$  be a schedule of  $Alg$  onto  $Arc$  with initial temperature  $T_{init}$ . If each task of  $Alg$  has been scheduled on the subset of cores  $Q_{best}$  defined by Eq. (19), thus satisfying the temperature constraint  $T_{obj}$ , and if  $T_{init} \leq T_{obj}$ , then the maximum temperature reached during one execution of  $S$  starting at  $T_{init}$  will also meet the  $T_{obj}$  constraint.

**Proof:** By hypothesis,  $T^{(0)} = T_{init} \leq T_{obj}$ . Then, the maximum temperature during  $S$  is equal to  $\max_{1 \leq i \leq n} T^{(i)}$ . Since each scheduling decision satisfies Eq. (19), it follows that  $\forall 1 \leq i \leq n, T^{(i)} \leq T_{obj}$ . As a conclusion we have  $\max_{1 \leq i \leq n} T^{(i)} \leq T_{obj}$ .  $\square$

#### D. Dealing with reactive systems

Propositions 3 and 4 are valid when the schedule is executed once, but also when the schedule is executed repeatedly and infinitely, as is the case for *reactive systems*. What characterizes a reactive system is that it controls some physical device (e.g., a satellite) and that it must continue to do so during the entire life of this physical device. Proposition 5 is valid when the schedule  $S$  is executed once, but not when it is repeated infinitely. The reason is due to the difference between the initial temperature  $T_{init}$  when the schedule starts and the final temperature  $T_f$  when the schedule ends (and also to the fact that the temperature curve depends on the initial temperature, as opposed to the GSFR and the power). Two cases arise:

- 1) If  $T_{init} < T_f \leq T_{obj}$ , then executing a second time the same schedule will inevitably *increase* further the temperature, so after some bounded number of executions of this schedule, the multicore temperature will violate the

$T_{obj}$  constraint. Recall that the cooling times are static and have been inserted in the schedule based on  $T_{init}$ . This is *not safe*.

- 2) If  $T_f < T_{init} \leq T_{obj}$ , then executing a second time the same schedule will inevitably *decrease* further the temperature, so after a large number of executions of this schedule, the multicore temperature will drop to the ambient temperature. This is *not optimal*.

Therefore, in order to be safe regarding the  $T_{obj}$  constraint and to be optimal, we should guarantee that  $T_f = T_{init}$ , which can only be achieved by being in Case 1 and then inserting on each core a cooling time until the average temperature of the multicore is equal to  $T_{init}$  (because we can cool down the multicore after executing the schedule by inserting a cooling time, while we cannot heat it). Proposition 6 generalizes Proposition 5 to the case of a schedule executed repeatedly.

**Proposition 6:** Let  $S$  be a schedule of  $Alg$  onto  $Arc$  with initial temperature  $T_{init}$ , final temperature  $T_f$ , and execution time  $C_{max}$ . If each task of  $Alg$  has been scheduled on the subset of cores  $Q_{best}$  defined by Eq. (19) (thus satisfying the temperature constraint  $T_{obj}$ ), if  $T_{init} \leq T_f \leq T_{obj}$ , and if we insert a cooling time of size  $\delta$  at the end of  $S$  such that  $T^{cool}(C_{max}, \delta) = T_{init}$ , then the maximum temperature reached during an arbitrary number of executions of  $S$  starting at  $T_{init}$  will also meet the  $T_{obj}$  constraint.

**Proof:** We prove this property by induction on the number  $m$  of executions of  $S$ . Let  $MaxTemp(k, S)$  denote the maximal temperature during the  $k$ -th execution of  $S$ .

The case  $MaxTemp(1, S) \leq T_{obj}$  is proved by Proposition 5. This first execution of  $S$  is followed by a cooling time of size  $\delta$ , hence  $T(C_{max} + \delta) = T_{init}$ , which is the start time of the second execution of  $S$ .

The induction hypothesis is then:

$$\max_{1 \leq k \leq m} MaxTemp(k, S) \leq T_{obj} \quad (21)$$

The  $m$ -th execution of  $S$  is followed by a cooling time of size  $\delta$ , hence  $T(m \cdot (C_{max} + \delta)) = T_{init}$ , which is the start time of the  $m+1$ -th execution of  $S$ . Applying the reasoning for  $MaxTemp(1, S)$  to the  $m+1$ -th execution yields  $MaxTemp(m+1, S) \leq T_{obj}$ . By the induction hypothesis, the proof is then concluded.  $\square$

The size  $\delta$  of the cooling time depends on the difference between  $T_f$  and  $T_{init}$ . It is obtained by solving for  $\delta$  the equation  $T^{cool}(C_{max}, \delta) = T_{init}$ :

$$\begin{aligned} \frac{B'}{A} + \left(T_f - \frac{B'}{A}\right) \cdot e^{-A(\delta - C_{max})} &= T_{init} \\ \iff e^{-A(\delta - C_{max})} &= \frac{T_{init} - B'/A}{T_f - B'/A} \\ \iff \delta &= C_{max} - \frac{1}{A} \cdot \log\left(\frac{T_{init} - B'/A}{T_f - B'/A}\right) \end{aligned}$$

Since  $\delta$  must be an integer number, we round it up:

$$\delta = \left\lceil C_{max} - \frac{1}{A} \cdot \log\left(\frac{T_{init} - B'/A}{T_f - B'/A}\right) \right\rceil \quad (22)$$

Finally, reactive systems must comply to hard deadlines. We do not directly address this when we generate the Pareto fronts. Once the Pareto front is computed, the user can eliminate all the points that fail to meet his or her hard deadline, and then choose one solution among the remaining ones by considering the other criteria.

### E. Taking into account the temperature of the adjacent cores

In a multicore, multiple cores are located on a single chip at a very short distance from each other, so the temperature of each core impacts the other cores. This is taken into account by Eqs. (14) and (15).

Now, one situation that can arise during our list scheduling algorithm is when the current task  $\tau^{(n)}$  is scheduled at step  $(n)$  on some core  $c$  such that  $c$ 's neighbors are (partly) idle during the duration of  $\tau^{(n)}$ . This is illustrated in Fig. 6(a) where task  $\tau^{(n)}$  is scheduled on  $c_2$ . The risk is that the temperature computed at the end of  $\tau^{(n)}$  is under-estimated because the tasks that will be scheduled on the neighbors of  $c_2$  (i.e.,  $c_1$  and  $c_3$  in Fig. 6(a)) in a *future* step of the heuristic will not be accounted for. For instance, Fig. 6(b) illustrates the case of a task  $\tau^{(n+1)}$  that is scheduled on  $c_1$  at step  $(n+1)$ , causing an increase of the temperature on  $c_2$  that was not taken into account when we scheduled  $\tau^{(n)}$  on  $c_2$ .

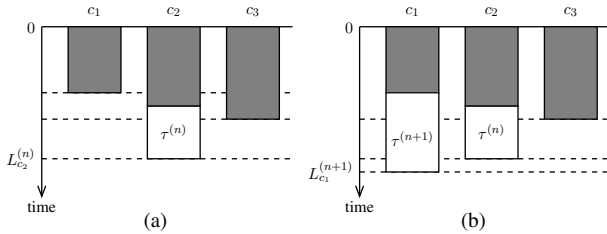


Fig. 6. (a) Partial schedule at step  $(n)$  and (b) at step  $(n+1)$ . A white box represents some new task  $\tau$  such that its vertical length is proportional to  $\mathcal{E}xe(\tau, c, f)$ . A gray box represents an arbitrary sequence of tasks scheduled during the previous steps.

We solve this issue by adding *virtual tasks* on all the neighbors to *over-estimate* the temperature: each time a task  $\tau^{(n)}$  is scheduled on some core  $c$ , for each neighbor  $c'$  of  $c$  such that the current finish time on  $c'$  is strictly less than the finish time on  $c$  (denoted  $L_c^{(n)}$  — note that it can be less than  $L_c^{(n)}$ ), we add on  $c'$  a virtual task that finishes exactly at  $L_c^{(n)}$  and that runs at frequency  $f_{max}$ . These virtual tasks modify the value of  $T_{c'}$  in Eq. (15), therefore guaranteeing that, whatever the future scheduling decisions, the runtime temperature on core  $c$  at time  $L_c^{(n)}$  will actually be below the temperature computed during the step  $(n)$  of our heuristic. This is illustrated in Fig 7. Of course, when actual tasks are scheduled on these cores  $c'$  during future steps, the virtual tasks are removed and the temperature is recomputed accordingly.

Table I summarizes the main computations used in ERPOT.

### F. Integer Linear Program

We now propose an ILP formulation of our scheduling problem, with the purpose of comparison with the heuristic algorithm presented in Section IV-B. The models and the

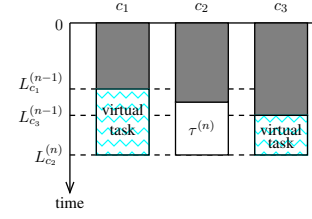


Fig. 7. Temperature over-estimation by adding virtual task to each of the neighbor of  $c_2$  because their respective finish time at step  $(n-1)$  was strictly less than  $L_{c_2}^{(n)}$ .

Execution time $\mathcal{E}xe(\tau, c, f) = \lceil \mathcal{E}xe_{nom}(\tau, c) / f \rceil$
Failure rate $\lambda_{sys} = \lambda_0 \cdot 10^{\frac{b(1-f)}{1-f_{min}}} \cdot e^{-\frac{E_a}{K} \left( \frac{1}{T(t)} - \frac{1}{T_0} \right)}$
Reliability $R(\tau, \mathcal{K}, t) = 1 - \left( \prod_{i=1}^k \left( 1 - e^{-\lambda_{sys}(c_i) \cdot \mathcal{E}xe(\tau, c_i, f_{c_i})} \right) \right)$ (computed with Reliability Block Diagrams)
GSFR $\Lambda(S) = -\log(R(S)) / U(S)$
Utilization $U(S) = \sum_{(\tau, c, f) \in S} \mathcal{E}xe(\tau, c, f)$
Power $P_{sys}(t) = \underbrace{\alpha \cdot T(t)}_{leakage} + \underbrace{\beta_h + \gamma \cdot C_{ef} \cdot V^2 \cdot f}_{dynamic}$
Temperature differential equation $C \cdot \left( \frac{dT_c(t)}{dt} \right) + G(T_c(t) - T_{amb}) + 2D_{heat} = P(t)$
Heat transfer from neighbor cores $2D_{heat} = \sum_{c' \in nbr(c)} \kappa(c, c') \cdot (T_c(t) - T_{c'}(t))$
Solution $T_c(t) = T_{\infty}^{heat} + (T_0 - T_{\infty}^{heat}) \cdot e^{-A(t-t_0)}$
Steady state temperature $T_{\infty}^{heat} = B/A$

TABLE I. SUMMARY OF ALL THE COMPUTATIONS.

assumptions used in Section IV-B are also used here for the ILP program. The decision variables are the following:

- $S_{ik} \in \mathbb{N}$  : start time of replica  $k$  of task  $i$
- $F_{ik} \in \mathbb{N}$  : finish time of replica  $k$  of task  $i$
- $Sb_{ik} \in \mathbb{N}$  : start time of replica  $k$  of data dependency  $i$
- $Fb_{ik} \in \mathbb{N}$  : finish time of replica  $k$  of data dependency  $i$
- $W \in \mathbb{N}$  : total execution time of the application

$$\begin{aligned}
 x_{ikc} &= \begin{cases} 1 & \text{if replica } k \text{ of task } i \text{ is assigned to core } c \\ 0 & \text{otherwise} \end{cases} \\
 x_{ikcfs} &= \begin{cases} 1 & \text{if replica } k \text{ of task } i \text{ is assigned to core } c \text{ at} \\ & \text{frequency } f \text{ and after a cooling time} \\ & \text{of } s \text{ time units} \\ 0 & \text{otherwise} \end{cases} \\
 \sigma_{ijkk'} &= \begin{cases} 1 & \text{if replica } k \text{ of task } i \text{ starts before replica } k' \\ & \text{of task } j \\ 0 & \text{otherwise} \end{cases} \\
 Y_{iK} &= \begin{cases} 1 & \text{if task } i \text{ is replicated } K \text{ times} \\ 0 & \text{otherwise} \end{cases} \\
 B_{ik} &= \begin{cases} 1 & \text{if replica } k \text{ of task } i \text{ has an outgoing} \\ & \text{data dependency} \\ 0 & \text{otherwise} \end{cases}
 \end{aligned}$$

The main objective of our optimization problem is minimizing the total execution time. Then, two kinds of ILP constraints

must be formulated. The first kind are the constraints that guarantee the schedulability:

- 1) Every replica  $k$  of task  $i$  should be assigned to exactly one core  $c$ :

$$\forall i, \forall k, \sum_c x_{ikcfs} = 1 \quad (23)$$

- 2) Every replica  $k$  of task  $i$  on core  $c$  should be assigned to exactly one level of frequency and be preceded by exactly one cooling time (possibly of size 0):

$$\forall i, \forall k, \forall c, \sum_{f,s} x_{ikcfs} = x_{ikc} \quad (24)$$

- 3) The finish time of every replica  $k$  of task  $i$  should be less or equal than the total execution time:

$$\forall i, \forall k, F_{ik} \leq W \quad (25)$$

- 4) The finish time of every replica  $k$  of task  $i$  is computed based on its execution time and its start time:

$$\forall i, \forall k, F_{ik} = S_{ik} + \sum_{c,f,s} x_{ikcfs} \cdot exe_c(i, c, f, s) + Fb_{ik} \quad (26)$$

where,  $exe_c(i, c, f, s)$  is the execution time of task  $i$  on core  $c$  at the  $f$ -th frequency level and after a cooling time of size  $s$ :  $exe_c(i, c, f, s) = \mathcal{E}xe(i, c, f) + s$ .

- 5) Tasks can not overlap and must obey their precedence order ( $M$  is a constant greater than the largest existing number in the ILP program — “big  $M$  method” [39]):

$$\forall i \neq j, \forall k, \forall k', \sigma_{ijkk'} + \sigma_{jik'k} \leq 1 \quad (27)$$

$$\forall i, \forall j, \forall k, \forall k', S_{ik} \leq S_{jk'} + (1 - \sigma_{ijkk'}) \cdot M \quad (28)$$

$$\forall i, \forall j, \forall k, \forall k', \forall c, F_{ik} \leq (2 - x_{ikc} - x_{jk'c}) \cdot M + S_{jk'} + (1 - \sigma_{ijkk'}) \cdot M \quad (29)$$

$$\forall i \in pred(j), \forall k, \forall k', F_{ik} \leq S_{jk'} \quad (30)$$

$$\forall i \in pred(j), \forall k, \forall k', \sigma_{ijkk'} = 1 \quad (31)$$

- 6) If task  $j$  is a successor of  $i$  and both are assigned to different cores, then this data dependency must be transmitted on the bus:

$$\forall i, \forall k, \forall j \in pred(i), \forall k', \forall c' \neq c, B_{ik} = \bigvee_c \left( x_{ikc} \wedge \left( \bigvee_{j,k',c'} x_{jk'c'} \right) \right) \quad (32)$$

where the logical operators  $\vee$  and  $\wedge$  are linearized [39].

- 7) The start time of data dependency  $i$  is computed based on the first idle time of the bus and on the previous data dependencies transmitted on the bus:

$$\forall i, \forall k, \forall b, Sb_{ik} = \sum_{j,k'} ((\sigma_{jik'k} \wedge B_{ik} \wedge B_{jk'}) \cdot exe_b(j, b)) \quad (33)$$

where  $exe_b(j, b)$  is the transmission time of data-dependency  $j$  on bus  $b$ :  $exe_b(j, b) = \mathcal{E}xe(j, b, f_b)$  (recall that buses operate at the fixed frequency  $f_b$ , and that we do not insert cooling times on the buses).

- 8) The finish time of each data dependency is the sum of its start time and its transmission time:

$$\forall i, \forall b, \forall k, Fb_{ik} = Sb_{ik} + B_{ik} \cdot exe_b(i, b) \quad (34)$$

- 9) Data dependencies must be serialized on the bus:

$$\forall i, \forall k, \forall j \geq i, \forall k', \forall b, Sb_{ik} \leq Sb_{jk'} - exe_b(i, b) + (1 - B_{ik} + \sigma_{ijkk'}) \cdot M \quad (35)$$

The second kind are the ILP constraints that guarantee that the GSFR / power consumption / temperature remain below  $\Lambda_{obj} / P_{obj} / T_{obj}$ :

- 1) The GSFR must be less than or equal to  $\Lambda_{obj}$ :

$$\forall i, \sum_k Y_{ik} = 1 \quad (36)$$

$$\forall i, \forall c, \sum_k x_{ikc} \leq 1 \quad (37)$$

$$\forall i, \sum_{k,c} x_{ikc} = \sum_k k \cdot Y_{ik} \quad (38)$$

$$\forall i, \sum_{k,c,f,s} x_{ikcfs} \cdot GSFR(c, f, s) + \sum_{k,b} B_{ik} \cdot GSFR(b, f_b, 0) \leq \Lambda_{obj} \quad (39)$$

- 2) The power consumption must be less than  $P_{obj}$ :

$$\sum_{i,k,c,f,s} x_{ikcfs} \cdot exe_c(i, c, f, s) \cdot P(f, s) + \sum_{i,k,b} B_{ik} \cdot P(f_b, 0) \cdot exe_b(i, b) \leq P_{obj} \cdot W \quad (40)$$

where  $P(f, s)$  is the sum of leakage and dynamic power consumption when the task runs at frequency  $f$  and is preceded by a cooling time of size  $s$ .

- 3) The temperature on each hardware component (cores and bus) must be less than or equal to  $T_{obj}$ :

$$\forall i, \forall k, \log(T_{\infty}^{heat} - T_0) - a \cdot F_{ik} + C \cdot M \geq \log(T_{\infty}^{heat} - T_{obj}) \quad (41)$$

$$\forall i, \forall k, \log(T_{\infty}^{cool} - T_0) - a \cdot F_{ik} \leq \log(T_{obj} - T_{\infty}^{cool}) + (1 - C) \cdot M \quad (42)$$

where  $T_0$ ,  $T_{\infty}^{heat}$ , and  $T_{\infty}^{cool}$  represent respectively the initial temperature at  $t_0$ , the heating steady state temperature, and the cooling steady state temperature. Eqs. (41) and (42) are for the cores; for the bus it suffices to replace  $F_{ik}$  by  $Fb_{ik}$  and to take the value of parameter  $a$  corresponding to the bus.

Based on these equations, the main objective of ILP is to minimize the total the execution length (the  $W$  variable in our ILP formulation), under the constraints specified by Eqs (23) to (42). In Section V-D, we will compare the Pareto fronts computed respectively by our quad-criteria heuristic ERPOT and by an ILP program.

## V. SIMULATION RESULTS

We ran several kinds of experiments to evaluate our ERPOT heuristic. In Section V-A, we assess the influence of the *temperature*, *power consumption*, and *reliability* constraints on the *execution time*. In Section V-B, we show a whole Pareto front for a given problem instance. In Section V-C, we compare ERPOT with the PowerPerf-PET scheduling heuristic from [4]. Finally, in Section V-D, we compare ERPOT with the ILP program of Section IV-F.

The target multicore chip is shown in Fig. 3(b) and the parameter values are provided in Table II, taken in part from [8] and [7].

$\lambda_0 = 10^{-5}$ , $C = 0.03 JK^{-1}$ , $G = 0.3 WK^{-1}$ , $\beta_h = -11 W$ , $\beta_c = -25 W$ , $\alpha = 0.1 WK^{-1}$	for each core
$C = 0.01 JK^{-1}$ , $G = 0.1 WK^{-1}$ , $\beta_h = -4 W$ , $\beta_c = -8 W$ , $\alpha = 0.04 WK^{-1}$	for the bus
$C_{ef} = 10^{-8} JV^{-2}$	same for the cores and the bus
$\kappa(\text{bus}, c_i) = 0.03 WK^{-1}$ , $\kappa(c_1, c_2) = \kappa(c_3, c_4) = 0.1 WK^{-1}$	thermal conductivity
$\{(900 MHz, 1.20 V), (600 MHz, 1.10 V), (300 MHz, 1.06 V)\}$	(voltage, frequency) pairs for the cores
$\{f_{max} = f_3 = 1, f_2 = \frac{2}{3}, f_{min} = f_1 = \frac{1}{3}\}$	scaling factors
$(300 MHz, 1.06 V)$	(voltage, frequency) pair for the bus
$f_b = \frac{1}{3}$	scaling factor

TABLE II. PARAMETER VALUES.

### A. Influence of the constraints on the schedules

Fig. 8 has been obtained with an  $\mathcal{Alg}$  graph consisting of 41 nodes, generated randomly with TGFF [40] (with the maximum value of in and out degree set to 4), and scheduled on the fully connected quad-core chip specified above. The nominal WCETs of the tasks are in the range  $[5 ms, 15 ms]$  while the nominal WCCTs of the data-dependencies are in the range  $[3 ms, 5 ms]^3$ .

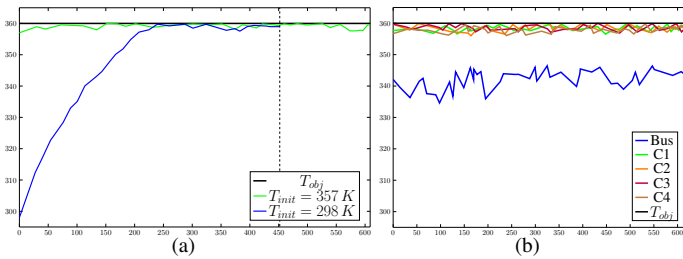


Fig. 8. (a) Evolution of the temperature when  $T_{init} = 298 K$  and  $T_{init} = 357 K$ . (b) Evolution of the temperature of each component.

Fig. 8(a) shows the variation of the chip temperature in function of the execution time and the effect of the insertion of cooling times in the schedule, for two different values of the initial temperature  $T_{init}$ : 298 K and 357 K. In both cases,  $T_{obj} = 360 K$ ,  $P_{obj} = 2 W$ , and  $\Lambda_{obj} = 10^{-8}$ . When  $T_{init} = 298 K$ , the temperature increases steadily during a transient phase, and then stabilizes just below  $T_{obj}$ , by virtue of the cooling times. When  $T_{init} = 357 K$ , the temperature

remains just below  $T_{obj}$  during the whole schedule, again by virtue of the cooling times. The initial temperature has a significant impact on the schedule length, from 451 ms for 298 K (indicated by the dashed vertical line) to 608 ms for 357 K, a 35% increase.

Fig. 8(b) depicts the temperature variation of the five hardware components of the chip (bus, C1, C2, C3, and C4) during a schedule produced with the same parameters as Fig. 8(a). The temperatures of the four cores remain in a very small interval,  $[356 K, 360 K]$ , demonstrating the effectiveness of our scheduling heuristic for the peak temperature. The bus temperature is significantly below for the simple reason that the bus is often idle. The fact that the temperature variations are very small, both over time and between the cores, is also very good to limit the aging of the chip [7].

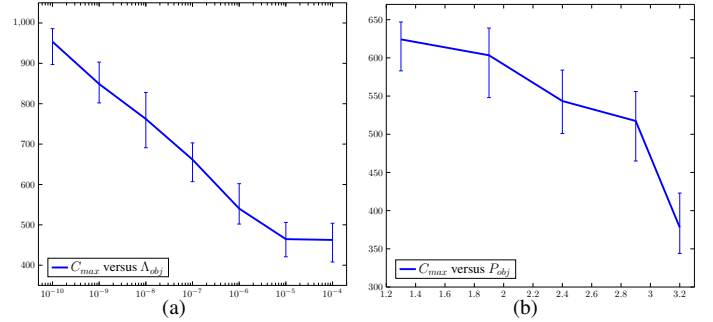


Fig. 9. (a) Influence of  $\Lambda_{obj}$  and (b) of  $P_{obj}$  on the execution time.

Fig. 9 has been obtained with 50 DAGs generated randomly, each with 50 tasks having an  $\mathcal{E}xe_{nom}$  in the range  $[3, 12]$ , and such that the total sum of the  $\mathcal{E}xe_{nom}$  of their tasks is in the range  $[540, 560]$ . Each point is the average value of the  $C_{max}$  over the 50 DAGs and each vertical bar shows the range around the average value. The schedule length increases when  $\Lambda_{obj}$  decreases (Fig. 9(a)). This is expected since more replications are required to satisfy the lower failure rate constraint: the two criteria are antagonistic. Moreover, the schedule length increases when  $P_{obj}$  decreases (Fig. 9(b)). This is expected since lowering the power consumption requires lowering the frequencies used by the cores, which increases the execution time. Again, the two criteria are antagonistic.

### B. Pareto fronts obtained with ERPOT

In this Section, we compute the whole Pareto front for an  $\mathcal{Alg}$  graph with 41 nodes onto the quad-core  $\mathcal{Arc}$  graph of Fig. 3(b) with the parameters of Table II. Ideally, we would like to visualize this Pareto front in 4D. However, when printed on paper, it is very hard to understand. To circumvent this difficulty, we show several Pareto fronts in 3D in the (execution time, GSFR, temperature) space and make it vary in the fourth dimension, the power consumption. We use 10 different values for each criterion. These threshold values must be provided by the user because they are application and platform dependent.

- $\Lambda_{obj} \in \{10^{-9}, 3.16 \cdot 10^{-9}, 10^{-8}, \dots, 3.16 \cdot 10^{-5}\}$ ;
- $P_{obj} \in \{1.3, 1.6, 1.9, \dots, 4.0\}$ , in Watts;

<sup>3</sup>From now on, the time unit will be the millisecond (ms).



- $T_{obj} \in \{340, 345, 350, \dots, 385\}$ , in Kelvin.

Algorithm 2 implements the grid method for our four criteria. The function ERPOT with the parameters  $\Lambda_i, P_j, T_k$  returns the Pareto point that minimizes the execution time under the constraints  $\Lambda < L_i$ ,  $P < P_j$ , and  $T < T_k$ . Since  $\Lambda_{obj}$  follows a logarithmic scale,  $\Lambda_{incr}$  is used as a multiplier.

---

**Algorithm 2** Grid method algorithm for 4 criteria.

---

**input:** The range  $[\Lambda_{min}, \Lambda_{max}]$  and the increment  $\Lambda_{incr}$   
**input:** The range  $[P_{min}, P_{max}]$  and the increment  $P_{incr}$   
**input:** The range  $[T_{min}, T_{max}]$  and the increment  $T_{incr}$   
**output:** The list of Pareto points  $Res$

- 1: **function** GRID( $\Lambda_{min}, \Lambda_{max}, \Lambda_{incr}, P_{min}, P_{max}, P_{incr}, T_{min}, T_{max}, T_{incr}$ )
- 2:  $Res \leftarrow \emptyset; \Lambda_1 \leftarrow \Lambda_{min}; i \leftarrow 1$
- 3: **while**  $\Lambda_i \leq \Lambda_{max}$  **do**
- 4:      $P_1 \leftarrow P_{min}; j \leftarrow 1$
- 5:     **while**  $P_j \leq P_{max}$  **do**
- 6:          $T_1 \leftarrow T_{min}; k \leftarrow 1$
- 7:         **while**  $T_k \leq T_{max}$  **do**
- 8:              $Res \leftarrow Res \cup \text{ERPOT}(\Lambda_i, P_j, T_k)$
- 9:              $T_k \leftarrow T_k + T_{incr}$
- 10:         **end while**
- 11:          $P_j \leftarrow P_j + P_{incr}$
- 12:     **end while**
- 13:      $\Lambda_i \leftarrow \Lambda_i \times \Lambda_{incr}$
- 14: **end while**
- 15: **return** REMOVE\_NONDOMINATED\_POINTS( $Res$ )
- 16: **end function**

---

Fig. 10 shows the resulting Pareto front in 3D for three different values of  $P_{obj}$ ,  $1.3W$ ,  $2.5W$ , and  $4.0W$ . A lower value of  $P_{obj}$  implies higher values for the  $C_{max}$ . This is expected because the power consumption and the execution time are antagonistic.

As expected also, when  $T_{obj}$  decreases, the  $C_{max}$  increases because more cooling times must be inserted and lower frequencies are chosen. For instance, in Fig. 10(a), at  $360K$  the  $C_{max}$  varies in the range  $[866ms, 1038ms]$ , while at  $340K$  it varies in the range  $[1041ms, 1222ms]$ .

When  $\Lambda_{obj}$  increases, the  $C_{max}$  increase because more tasks must be replicated to compensate for the higher failure rate. Again this is expected because these two criteria are antagonistic. For instance, in Fig. 10(c), at  $10^{-5}$  failures per  $ms$ , the  $C_{max}$  varies in the range  $[341ms, 498ms]$ , while it varies in  $[403ms, 594ms]$  at  $10^{-9}$ .

### C. Comparison with PowerPerf-PET

We have also compared ERPOT with the PowerPerf-PET heuristic from [4] (Algorithm 7) but without considering the reliability since PowerPerf-PET does not address this criterion. PowerPerf-PET uses two separate cost functions to select the core and the frequency to execute the current task. To select the core, it evaluates, for each task  $\tau_i$ , the product of the total power consumption of each core before mapping  $\tau_i$ , and its earliest possible available time for executing  $\tau_i$ . The task is allocated to the core having the minimum value of this product (the ‘‘PowerPerf’’ part). Then, to select the frequency,

it uses a weighted sum of the performance  $P$ , the energy  $E$ , and the temperature  $T$  (the ‘‘PET’’ part). In contrast, we use a unique cost function to select the core, its frequency, and the length of the cooling time (if any). Regarding the cooling times, PowerPerf-PET never inserts one. Moreover, the temperature model of PowerPerf-PET is based on measurement rather than an analytic model based on the differential heat propagation equation, and it does not take into account the heat propagation from the neighbor cores. Similarly, the power consumption model is based on measurement, and the effect of the temperature on the power consumption is not taken into account.

As application graphs, we choose the five benchmarks from the E3S suite [41] and five DAGs randomly generated with TGFF [40] (with the maximum value of in and out degree set to 4). For each DAG, the target architecture is the quad-core platform of Section V-A. For ERPOT, we take the following values for the  $P_{obj}$  and  $T_{obj}$  constraints, resulting in 100 points in each Pareto front:

- $P_{obj} \in \{1.3, 1.6, 1.9, \dots, 4.0\}$ , in Watts;
- $T_{obj} \in \{340, 345, 350, \dots, 385\}$ , in Kelvin.

For PowerPerf-PET, the three weights of the ‘‘PET’’ weighted sum are each taken in the interval  $[0, 1]$  with a 0.01 increment.

Thanks to the grid method, ERPOT produces one Pareto point in each cell of this 2D space. This is not the case of PowerPerf-PET because it relies on the transformation method with a weighted sum. As explained in Section II, this does not allow exploring the entire search space. Table III reports the number of cells for which each algorithm succeeds in finding a valid schedule. In each cell of the grid containing a solution from ERPOT and from PowerPerf-PET, we compute the percentage between the  $C_{max}$  of these two schedules as:

$$\frac{C_{max}(\text{PowerPerf-PET}) - C_{max}(\text{ERPOT})}{C_{max}(\text{PowerPerf-PET})} \times 100 \quad (43)$$

Finally, we compute the average of these percentages over all the suitable cells of the 2D space. The results are reported in Table III. ERPOT systematically outperforms PowerPerf-PET by at least 27%. Several reasons explain this significant difference. First, PowerPerf-PET is based on a weighted sum of its three criteria  $P$ ,  $E$ , and  $T$ . This does not allow the concave parts of the Pareto front to be found (see Fig. 1(b)). As a consequence, PowerPerf-PET computes the convex hull of the Pareto front while ERPOT computes the actual one, including its concave parts. Second, ERPOT uses a smart cost function to sort the ready tasks, taking into account for each task  $\tau$  the longest execution path from  $\tau$  to the end of the graph (see Section IV-B). This allows us to schedule first the tasks that are in the critical path, which reduces the overall execution time.

### D. Evaluation of the ILP model

We have implemented our ILP program (see Section IV-F) in the CPLEX ILOG solver [42] – version 12.6.3, and we have run it on an Intel quad-core i5 CPU with 6GB RAM. It returns the optimal result, i.e., the schedule with the *minimal*

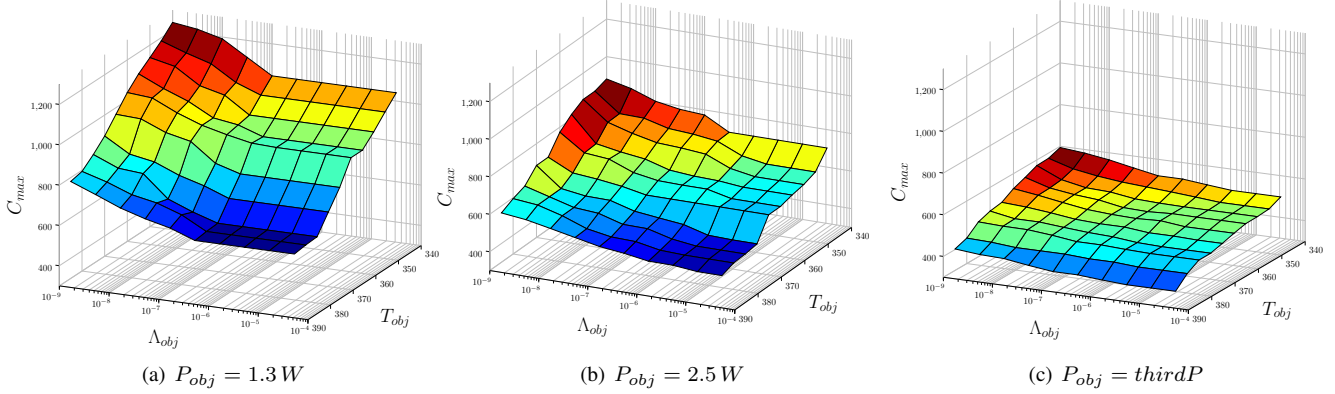


Fig. 10. Pareto fronts in 3D for three different values of  $P_{obj}$ :  $1.3W$ ,  $2.5W$ , and  $4.0W$ .

Benchmark	automotive	consumer	networking	office	telecom	random	random	random	random	random
DAG size	24	12	13	5	30	40	50	60	70	80
ERPOT (cells)	100	100	100	100	100	100	100	100	100	100
PowerPerf-PET (cells)	7	4	6	3	8	7	8	6	5	7
$C_{max}$ improvement (%)	<b>27.64</b>	36.88	38.30	<b>39.92</b>	29.33	30.81	35.82	34.59	31.32	30.67

TABLE III. ERPOT VS. POWERPERF-PET: ERPOT SYSTEMATICALLY OUTPERFORMS POWERPERF-PET ON THE  $C_{max}$  (BY AN AVERAGE OF 33.5%).

execution time under the  $\Lambda_{obj}$ ,  $P_{obj}$ , and  $T_{obj}$  constraints. The drawback is that the complexity of finding this optimal schedule is exponential in the size of the problem instance (number of tasks of the  $Alg$  graph plus number of cores times number of frequencies). To be specific, our ILP program was not able to complete its execution for DAGs larger than 9 tasks because the CPLEX solver ran out of memory.

We have run our ILP program on 10 DAGs randomly generated with TGFF [40], each with 8 tasks, and a homogeneous dual-core with a single bus with three frequency/voltage levels. The WCETs of the tasks are randomly chosen in the range  $[3ms, 12ms]$  while the WCCTs are randomly chosen in the range  $[2ms, 4ms]$ . Besides, the cooling times are limited to  $1ms$ . Finally, ten different values of each criterion  $\Lambda_{obj}$ ,  $P_{obj}$ , and  $T_{obj}$  are considered (as in Section V-B).

For each DAG, we build the full 4D Pareto front with ERPOT and with the ILP program, and in each cell we compute the percentage between the  $C_{max}$  of these two schedules as:

$$\frac{C_{max}(\text{ERPOT}) - C_{max}(\text{ILP})}{C_{max}(\text{ERPOT})} \times 100 \quad (44)$$

For each Pareto front, we compute the minimum, maximum, and average difference between the two solutions. Table IV summarizes the results. On average, the length of the non-optimal schedule obtained with our ERPOT heuristic is between 8% and 10% above the length of the optimal schedule obtained with the ILP program, which we claim is not too bad. However, recall that the ILP solver can only compute the Pareto front for very small DAGs, no larger than 8 tasks.

Finally, Figure 11 plots the percentage of the  $C_{max}$  between the two Pareto fronts generated by ERPOT and by the ILP program for the DAG # 6 from Table IV. The largest deviations between ERPOT and ILP occur when the  $\Lambda_{obj}$ ,  $P_{obj}$ , and  $T_{obj}$

DAG	1	2	3	4	5	6	7	8	9	10
min (%)	3.20	2.90	3.12	3.08	2.89	3.84	2.71	4.16	3.94	3.30
max (%)	23.74	24.10	25.60	22.38	24.67	25.10	25.47	23.74	22.69	26.38
avg. (%)	8.38	9.26	<b>8.17</b>	8.55	8.49	9.43	9.07	<b>9.94</b>	8.92	9.33

TABLE IV. ERPOT VS. ILP: ILP SYSTEMATICALLY OUTPERFORMS ERPOT ON THE  $C_{max}$  (BY AN AVERAGE OF 8.95%).

constraints are the more stringent. The reason is that the ILP program makes better choices between inserting cooling times and lowering the frequency/voltage.

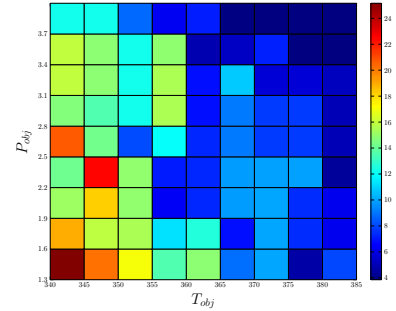


Fig. 11. Heatmap of the percentage between the  $C_{max}$  obtained by ERPOT and by the ILP program ( $\Lambda_{obj} = 10^{-6}$ ).

## VI. RELATED WORK

Several related works optimize the executing time, reliability, power consumption, and temperature for applications running on multicores. Most of them consider only two criteria, the execution time and one of the other three criteria. A

few related work consider three criteria. However, no existing results consider the four criteria altogether.

Many results address the problem in the context of applications modeled as a set of real-time tasks, usually pre-emptible, and scheduled by a real-time operating system (RTOS) according to some priority policy (see e.g. [43]–[45] to cite only a few). Each task  $\tau_i$  is defined by a tuple  $(A_i, C_i, D_i, \pi_i)$ , where  $A_i$  is the arrival time (defined either according to a periodic or sporadic activation model),  $C_i$  is the worst-case execution time,  $D_i$  is the deadline, and  $\pi_i$  is the priority. Since our application model is totally different, we do not detail these works.

In [5], Girault and Kalla present a bi-criteria optimization ready list heuristic algorithm to schedule a DAG of tasks onto a heterogeneous multi-core processor. The algorithm minimizes both the total execution time and the soft error rate. Instead of directly using the system’s reliability as an optimization criterion, the authors introduce a new criterion called the Global System Failure Rate (GSFR). The GSFR is computed based on the system’s reliability (i.e., the reliability of the schedule on the multicore, computed with classical reliability techniques such as reliability block diagrams) and on the total execution time. The main advantage of the GSFR over the reliability is that it is an *invariant* measure of the schedule (see Section III-C for details), which makes it suitable to use the transformation method. This allows the authors to use the transformation method to compute the Pareto front in the 2D space (execution time, GSFR). This method has been generalized in [2] by Assayad et al. to take into account the power consumption, therefore providing a tri-criteria list scheduling algorithm to optimize the execution time, the GSFR, and the power consumption. The effect of voltage and frequency on the failure rate per time unit of the cores is taken into account. However, it does not take into account the temperature of the cores. ERPOT extends [2] precisely to take into account the temperature.

In [7], Das et al. propose a bi-criteria genetic algorithm to schedule a DAG of tasks onto a set of identical cores interconnected in a mesh network topology. The algorithm maximizes two criteria, (i) the system reliability (called the “performability”) and (ii) the lifetime, under a given energy constraint  $E_{max}$  and a given latency constraint  $P_{max}$  (the latency is incorrectly denoted “period”). DVFS is used to lower the total energy consumed. The reliability model, a variant of the Poisson model of Shatz and Wang [24], takes into account the voltage/frequency effect on the failure rate, as in [10]. The lifetime is computed by taking into account failures due to electromigration (EM), with a Weibull distribution. In order to improve the system reliability and the lifetime, some tasks are chosen to be actively replicated; this choice is made by the genetic algorithm. The result is a set of non-dominated solutions in the 2D space (reliability, lifetime). The reliability is improved by increasing the number of replicas, while the lifetime is improved by lowering the temperature. It follows that increasing the number of replicas increases the chip temperature, which in turn decreases the lifetime; in this sense the two criteria are antagonistic. However, due to the very high cost of the genetic algorithm, only small DAGs can be scheduled (up to 20 tasks), while we are able to handle DAGS

of size greater than 100 tasks. Besides, the effect of the chip temperature on the lifetime and on the reliability is not taken into account. Finally, the leakage power is ignored.

In [4], Sheikh and Ahmad address the PETOS problem (Performance, Energy, and Temperature Optimized Scheduling), where a DAG of tasks must be scheduled onto a set of  $M$  parallel cores operating under  $K$  available frequency levels. Because large DAGs are considered, only heuristic algorithms can be used (i.e., neither ILP nor exhaustive search algorithms). The authors propose 16 different heuristic, which are classified according to (i) the core selection strategy and (ii) the frequency selection strategy. None of the heuristic algorithms proposed in [4] is able to optimize the reliability, including PowerPerf-PET already described in Section V-C.

In [9], Xie et al. present an energy-efficient fault-tolerant list scheduling heuristic. The application model is a DAG of tasks, and the architecture model is a distributed memory multi-processor with a CAN network. Processors are heterogeneous and equipped with DVFS, but leakage power is ignored. The reliability model is the Poisson model of Shatz and Wang [24], and the effect of DVFS on the failure rate per time unit is taken into account as in [10], but not the temperature. Active replication is applied to each task of the DAG so as to satisfy a given reliability goal for the resulting system (e.g., 0.99). When doing so, the frequencies of the processors the task is mapped to are taken into account. The proposed heuristic minimizes the total energy under this reliability constraint, but the authors do not compute Pareto fronts.

Finally, Papadimitriou and Yannakakis address the problem of computing an *approximate* Pareto front in the case where the size of the exact Pareto front (i.e., the number of Pareto optima) is exponential in the size of the problem instance. They define the notion of  $\varepsilon$ -approximated Pareto front, for which each point is at most at a distance  $\varepsilon$  from an optimal Pareto point in each dimension (using the  $L_\infty$  norm). For this reason, the size of the  $\varepsilon$ -approximate Pareto front is much smaller than that of the exact Pareto front. The authors prove that for any  $n$ -criteria optimization problem, any problem instance  $x$ , and any  $\varepsilon$ , there exists an approximate  $\varepsilon$ -approximate Pareto front the size of which is polynomial in the size of  $x$  and in  $1/\varepsilon$ , but exponential in  $n$ .

## VII. CONCLUSION

We have presented a novel quad-criteria distributed scheduling heuristic called ERPOT (for Execution time, Reliability, Power consumption and Temperature), which minimizes the schedule length under three constraints: the power consumption, the maximal temperature, and the Global System Failure Rate (GSFR, which generalizes the classical failure rate per time unit of hardware elements to a whole schedule on a multicore architecture). These four criteria are all crucial to optimize embedded systems. By varying the three constraints and repeatedly invoking our ERPOT heuristic, we are able to compute the whole Pareto front in the 4D space (execution time, failure rate, power, temperature).

Using ERPOT in practice involves (i) modelling the application as a DAG of tasks, (ii) evaluating the WCET of each task

with a dedicated tool, (iii) gathering all the parameter values from the chip, (iv) building the Pareto front, and (v) choosing one solution from the Pareto front according to the application and user constraints.

The failure rate constraint is met by adding active replica in the schedule. Hence the failure rate and the schedule length are antagonistic criteria. The power consumption constraint is met by using Dynamic Voltage and Frequency Scaling (DVFS). Hence the schedule length and the power consumption are antagonistic criteria. Finally, the temperature constraint is met by inserting cooling times in the schedule (but also by lowering the voltage). Hence the schedule length and the temperature are antagonistic criteria.

The antagonisms between the criteria already make the scheduling problem very complex. Moreover, there are other interplays that must also be taken into account. For instance, lowering the voltage makes the hardware sensitive to lower energy particles, thereby increasing the nominal failure rates of the hardware components of the target architecture. ERPOT is the first scheduling heuristic able to take into account all those antagonisms.

Extensive experimental results show that our scheduling heuristic works very well: (i) on small application graphs, ERPOT is outperformed on average by less than 10% by an ILP program that produces the *optimal* Pareto fronts; (ii) on large application graphs, both synthetic and real-life, ERPOT outperforms the PowerPerf-PET scheduling heuristic on average by 33%. The largest deviations between ERPOT and ILP occur when the  $\Lambda_{obj}$ ,  $P_{obj}$ , and  $T_{obj}$  constraints are more stringent. The reason is that the ILP program makes better choices between inserting cooling times and lowering the frequency/voltage. This hints at potential avenues for future improvements of ERPOT.

It is tempting to extend our method to a general N-constraints method. However, in the context of real-time embedded systems, we believe that it is not possible. First, the constraints are inter-dependent, as evidenced by Eq. (6). The only possibility to get a general N-constraint method would be if the constraints were independent of each other. Second, incorporating the power consumption into our method required inserting virtual tasks in the schedule to fill holes, in order to avoid under-estimating the power consumption (see Sec. IV-C). Third, a similar issue emerged when incorporating the temperature, which required us to also insert virtual tasks to avoid under-estimating the peak temperature (see Sec. IV-E). Although it may seem that the solution is identical for the power consumption and temperature, this is not the case. As a matter of fact, keeping the system under a temperature threshold also requires inserting cooling times in the schedule.

## REFERENCES

- [1] L. Torres, P. Benoit, G. Sassatelli, M. Robert, F. Clermidy, and D. Puschini, "An introduction to multi-core system on chip – trends and challenges," in *Multiprocessor System-on-Chip*, pp. 1–21, Springer Science Business Media, Nov. 2010.
- [2] I. Assayad, A. Girault, and H. Kalla, "Tradeoff exploration between reliability, power consumption, and execution time for embedded systems," *International Journal on Software Tools for Technology Transfer*, vol. 15, no. 3, pp. 229–245, 2013.
- [3] Joint Electron Device Engineering Council, "Failure mechanisms and models for semiconductor devices," Tech. Report JEP 122-H, JEDEC, Aug. 2016.
- [4] H. Sheikh and I. Ahmad, "Sixteen heuristics for joint optimization of performance, energy, and temperature in allocating tasks to multi-cores," *ACM Trans. on Parallel Computing*, vol. 3, pp. 1–29, Aug. 2016.
- [5] A. Girault and H. Kalla, "A novel bicriteria scheduling heuristics providing a guaranteed global system failure rate," *IEEE Trans. on Dependable and Secure Computing*, vol. 6, pp. 241–254, Oct. 2009.
- [6] R. Viswanath, V. Wakharkar, A. Watwe, and V. Lebonheur, "Thermal performance challenges from silicon to systems," *Intel Technology Journal*, vol. Q3, 2000.
- [7] A. Das, A. Kumar, B. Veeravalli, C. Bolchini, and A. Miele, "Combined DVFS and mapping exploration for lifetime and soft-error susceptibility improvement in MPSoCs," in *Design, Automation & Test in Europe, DATE'14*, (Dresden, Germany), Mar. 2014.
- [8] P. Kumar and L. Thiele, "Thermally optimal stop-go scheduling of task graphs with real-time constraints," in *Asia and South Pacific Design Automation Conference, ASP-DAC'11*, IEEE, Jan. 2011.
- [9] G. Xie, Y. Chen, and X. Xiao, "Energy-efficient fault-tolerant scheduling of reliable parallel applications on heterogeneous distributed embedded systems," *IEEE Trans. on Sustainable Computing*, June 2017.
- [10] D. Zhu, R. G. Melhem, and D. Mossé, "The effects of energy management on reliability in real-time embedded systems," in *ICCAD'04*, pp. 35–40, IEEE / ACM, 2004.
- [11] M. Andjelkovic, M. Krstic, R. Kraemer, V. Veeravalli, and A. Steininger, "A critical charge model for estimating the SET and SEU sensitivity: A muller C-element case study," in *Asian Test Symposium, ATS'17*, pp. 82–87, IEEE, Nov. 2017.
- [12] Y. Haimes, L. Lasdon, and D. Wismer, "On a bicriterion formulation of the problems of integrated system identification and system optimization," *IEEE Trans. Systems, Man, and Cybernetics*, vol. 1, pp. 296–297, 1971.
- [13] M. Laumanns, L. Thiele, and E. Zitzler, "An efficient, adaptive parameter variation scheme for metaheuristics based on the epsilon-constraint method," *European J. of Operational Research*, vol. 169, no. 3, pp. 932–942, 2006.
- [14] Y. Xie and W.-L. Hung, "Temperature-aware task allocation and scheduling for embedded multiprocessor systems-on-chip (MPSoC) design," *The Journal of VLSI Signal Processing*, vol. 45, pp. 177–189, Dec. 2006.
- [15] M. Garey and D. Johnson, *Computers and Intractability, a Guide to the Theory of NP-Completeness*. San Francisco: W.H. Freeman Company, 1979.
- [16] V. T'Kindt and J. Billaut, *Multicriteria Scheduling – Theory, Models and Algorithms*. Springer, 2006.
- [17] I. Das and J. Dennis, "Normal-boundary intersection: A new method for generating the Pareto surface in nonlinear multicriteria optimization problems," *SIAM J. Opt.*, vol. 8, pp. 631–657, Mar. 1998.
- [18] P. Cousot and R. Cousot, "Abstract interpretation: A unified lattice model for static analysis of programs by construction or approximation of fixpoints," in *Principles of Programming Languages, POPL'77*, (Los Angeles (CA), USA), ACM SIGPLAN, Jan. 1977.
- [19] A. Colin and I. Puaut, "Worst case execution time analysis for a processor with branch prediction," *Real-Time Systems*, vol. 18, no. 2/3, pp. 249–274, 2000.
- [20] H. Theiling, C. Ferdinand, and R. Wilhelm, "Fast and precise WCET prediction by separate cache and path analyses," *Real-Time Systems*, vol. 18, pp. 157–179, May 2000.
- [21] S. Altmeyer, R. Davis, L. Indrusiak, C. Maiza, V. Nélis, and J. Reineke, "A generic and compositional framework for multicore response time analysis," in *International Conference on Real Time Networks and Systems, RTNS'15*, (Lille, France), pp. 129–138, ACM, Nov. 2015.
- [22] H. Rihani, M. Moy, C. Maiza, R. Davis, and S. Altmeyer, "Response time analysis of synchronous data flow programs on a many-core

- processor,” in *International Conference on Real-Time Networks and Systems, RTNS'16*, (Brest, France), pp. 67–76, ACM, Oct. 2016.
- [23] R. Davis, S. Altmeyer, L. Indrusiak, C. Maiza, V. Nelis, and J. Reineke, “An extensible framework for multicore response time analysis,” *Real-Time Syst.*, vol. 54, no. 3, pp. 607–661, 2018.
- [24] S. Shatz and J.-P. Wang, “Models and algorithms for reliability-oriented task-allocation in redundant distributed-computer systems,” *IEEE Trans. Reliability*, vol. 38, pp. 16–26, Apr. 1989.
- [25] A. Avizienis, J.-C. Laprie, B. Randell, and C. Landwehr, “Basic concepts and taxonomy of dependable and secure computing,” *IEEE Trans. on Dependable and Secure Computing*, vol. 1, pp. 11–33, Jan. 2004.
- [26] H. Balaban, “Some effects of redundancy on system reliability,” in *National Symposium on Reliability and Quality Control*, (Washington (DC), USA), pp. 385–402, Jan. 1960.
- [27] D. Rossi, M. Omana, C. Metra, and A. Paccagnella, “Impact of bias temperature instability on soft error susceptibility,” *IEEE Trans. Very Large Scale Integration Systems*, vol. 23, pp. 743–751, apr 2015.
- [28] A. M. Fard, M. Ghasemi, and M. Kargahi, “Response-time minimization in soft real-time systems with temperature-affected reliability constraint,” in *2015 CSI Symposium on Real-Time and Embedded Systems and Technologies, RTEST'15*, IEEE, oct 2015.
- [29] S. Hsueh, R. Huang, and C. Wen, “TASSER: A temperature-aware statistical soft-error-rate analysis framework for combinational circuits,” in *Fifteenth International Symposium on Quality Electronic Design*, IEEE, mar 2014.
- [30] J. Srinivasan, S. V. Adve, P. Bose, and J. A. Rivers, “Exploiting structural duplication for lifetime reliability enhancement,” in *ISCA*, pp. 520–531, IEEE, 2005.
- [31] M. Moy, C. Helmstetter, T. Bouhadiba, and F. Maraninchi, “Modeling power consumption and temperature in TLM models,” *Leibniz T. on Embedded Systems*, vol. 3, no. 1, pp. 1–29, 2016.
- [32] F. Kreith, *CRC Handbook of Thermal Engineering*. Mechanical and Aerospace Engineering Series, CRC Press, 1999.
- [33] T. Chantem, R. P. Dick, and X. S. Hu, “Temperature-aware scheduling and assignment for hard real-time applications on MPSoCs,” *IEEE Trans. Very Large Scale Integration Systems*, vol. 19, no. 10, pp. 1884–1897, 2011.
- [34] J. Knight and N. Leveson, “An experimental evaluation of the assumption of independence in multi-version programming,” *IEEE Trans. Software Engin.*, vol. 12, no. 1, pp. 96–109, 1986.
- [35] P. Jensen and M. Bellmore, “An algorithm to determine the reliability of a complex system,” *IEEE Trans. Reliability*, vol. 18, pp. 169–174, Nov. 1969.
- [36] A. S. Hartman, D. E. Thomas, and B. H. Meyer, “A case for lifetime-aware task mapping in embedded chip multiprocessors,” in *Proceedings of the eighth IEEE/ACM/IFIP international conference on Hardware/software codesign and system synthesis, CODES/ISSS'10*, ACM, 2010.
- [37] Z. Wang, S. Ranka, and P. Mishra, “Efficient task partitioning and scheduling for thermal management in multicore processors,” in *ISQED'15*, IEEE, 2015.
- [38] Y.-K. Kwok and I. Ahmad, “Static scheduling algorithms for allocating directed task graphs to multiprocessors,” *ACM Computing Surveys*, vol. 31, no. 4, pp. 406–471, 1999.
- [39] B. McCarl and T. Spreen, *Applied Mathematical Programming Using Algebraic Systems*. College Station (TX), USA: Texas A&M University, 2007.
- [40] “Task graphs for free.” <http://ziyang.eecs.umich.edu/~dickrpt/gtff>. Accessed: 2016-10-22.
- [41] “Embedded microprocessor benchmark consortium.” <http://www.eembc.org>.
- [42] IBM, “ILOG CPLEX optimizer.” <https://www-01.ibm.com/software/commerce/optimization/cplex-optimizer>. Accessed: 2016-10-22.
- [43] L. Huang, F. Yuan, and Q. Xu, “Lifetime reliability-aware task allocation and scheduling for MPSoC platforms,” in *Design Automation and Test in Europe Conference, DATE'09*, (Nice, France), pp. 51–56, Mar. 2009.
- [44] X. Qin, W. Wang, and P. Mishra, “TCEC: Temperature and energy-constrained scheduling in real-time multitasking systems,” *IEEE Trans. Computer-Aided Design of Integrated Circuits and Systems*, vol. 31, pp. 1159–1168, Aug. 2012.
- [45] Y. Ma, T. Chantem, and R. P. Dick, “Improving system-level lifetime reliability of multicore soft real-time systems,” *IEEE Trans. Very Large Scale Integration Systems*, vol. 25, pp. 1895–1905, Mar. 2017.

PLACE  
PHOTO  
HERE

**Athena Abdi** Athena Abdi received the B.Sc. and M.Sc. degrees in computer engineering from Amirkabir University of Technology (Tehran Polytechnic), Tehran, Iran in 2010 and 2012 respectively. Currently, she is a PhD candidate at Design and Analysis of Dependable Systems (DADS) laboratory of Amirkabir University of Technology. Her research interests include embedded and real-time systems, scheduling, optimization and fault tolerant design.

PLACE  
PHOTO  
HERE

**Alain Girault** Alain Girault holds a senior research fellow position at INRIA, the French National Institute for Research in Computer Science and Control. He received his PhD from the National Polytechnic Institute of Grenoble in 1994 after doing his research in the Verimag laboratory. He spent two years and a half as a postdoc researcher, in the ESTEREL team at INRIA in Sophia-Antipolis, France, in the PTOLEMY group at UC Berkeley, and in the PATH project at UC Berkeley. His research interests include the design of reactive systems, with a special concern

for distributed implementation, fault-tolerance, reliability, and low power.

PLACE  
PHOTO  
HERE

**Hamid Zarandi** Hamid R. Zarandi received his B.Sc., M.Sc., and Ph.D. degrees all in the department of computer engineering at Sharif University of Technology (SUT), Tehran, Iran, in 2000, 2002, and 2007, respectively. He is currently an associate professor in computer engineering and information technology department at Amirkabir University of Technology (AUT) (Tehran Polytechnic), since 2007. His research interests include dependability evaluation using fault injection techniques, fault-tolerant computing, dependable computer architecture, high

performance computing, and fault-tolerant embedded and real-time systems, on which he has published more than 100 referred conference and journal papers. Dr. Zarandi established the “Design and Analysis of Dependable Systems (DADS)” laboratory at Amirkabir University in 2007, and has chaired the laboratory since then. He is a member of the IEEE Computer Society, and the Computer Society of Iran (CSI).