

Temperature-Aware Scheduling and Assignment for Hard Real-Time Applications on MPSoCs

Thidapat Chantem, *Student Member, IEEE*, X. Sharon Hu, *Senior Member, IEEE*
and Robert P. Dick, *Member, IEEE*

Abstract—Increasing integrated circuit (IC) power densities and temperatures may hamper multiprocessor system-on-chip (MPSoC) use in hard real-time systems. This article formalizes the temperature-aware real-time MPSoC assignment and scheduling problem and presents an optimal phased steady-state mixed integer linear programming based solution that considers the impact of scheduling and assignment decisions on MPSoC thermal profiles to directly minimize the chip peak temperature. We also introduce a flexible heuristic framework for task assignment and scheduling that permits system designers to trade off accuracy for running time when solving large problem instances. Finally, for task sets with sufficient slack, we show that inserting idle times between task executions can further reduce the peak temperature of the MPSoC quite significantly.

Index Terms—Hard real-time systems, multiprocessor-system-on-chips, MPSoCs, scheduling and assignment, temperature-aware system-level design

I. INTRODUCTION

MULTIPROCESSOR systems-on-chips (MPSoCs) are now widely used in application-specific systems and high-performance computing. They offer performance, power consumption, and implementation complexity advantages over highly superscalar uniprocessor architectures. Their use, and scale, will increase dramatically in the coming years. According to Milchman [24], 16-core processors will be common within the next four years. Intel plans to deliver processors that have dozens or hundreds of cores during the next decade [3]. The use of heterogeneous MPSoCs can sometimes dramatically improve performance and power consumption relative to homogenous MPSoCs [16]. However, it can also increase complexity. It is likely that some future MPSoCs will be homogeneous and some will be heterogeneous. With the current use of MPSoCs in soft real-time applications such as gaming [44], it is expected that many hard real-time applications will soon be implemented using MPSoCs. In fact, FreeScale is now offering the QorIQ Embedded Multicore

Processor [29] that is intended for application domains that require real-time computing, e.g., aerospace applications.

MPSoC temperature is a strong function of power density. Increasing transistor counts and aggressive frequency scaling result in a significant increase in chip power density and temperature. Increasing chip temperature has significant impact on other design metrics including reliability, performance, and cost, as microprocessor failure rate depends exponentially upon operating temperature [39]. A 10–15 °C difference in operating temperature can result in a 2× difference in the lifespan of a device [42]. Temperature also affects speed; reduction of charge carrier mobility in transistors and increased interconnect latency resulting from high temperature degrade performance, requiring reduced clock frequencies or, worse yet, resulting in run-time failures.

Increasing power densities make package and cooling design for the worst-case scenario prohibitively expensive, since the cost of cooling solutions increases super-linearly with power consumption [12]. It is therefore necessary to design packaging and cooling solutions based on less than worst-case thermal profiles and compensate by preventing, hopefully rare, dangerous thermal scenarios at run-time. Most popular approaches react to critical temperatures by reducing frequency and voltage (i.e., performing hardware throttling), or by temporarily preventing instruction issue to reduce the power consumption, and hence temperature, of the processor [19].

Since the execution times of real-time tasks, and hence total system utilization, tend to vary significantly due to factors such as conditional branches and system inputs [46], real-time applications can exhibit great temperature variation at run-time. When the system utilization is low, the MPSoC may not have a high temperature problem, thanks to the amount of slack available in the system. On the other hand, a system with high utilization can push an MPSoC to its thermal limit [5], [43], [45]. In the worst case, the host MPSoC may lack run-time thermal management, leading to overheating and signal timing violations or permanent failure. More subtly, even those real-time systems containing MPSoCs that support run-time thermal management may fail when a temperature bound is reached, but for a different reason.

Most run-time thermal management techniques use thermal sensors to detect when the maximum safe temperature is approached and react by decreasing processor power consumption, e.g., by decreasing frequency or stalling instruction issue. These techniques share a common weakness: they decrease performance. If a real-time task running on an MPSoC with

Manuscript received August 2, 2009; accepted June 15, 2010.

T. Chantem is with the department of Computer Science & Engineering at the University of Notre Dame, Notre Dame, IN 46556, tchantem@nd.edu

X.S. Hu is with the department of Computer Science & Engineering at the University of Notre Dame, Notre Dame, IN 46556, shu@nd.edu

R.P. Dick is with the department of Electrical Engineering and Computer Science at the University of Michigan, Ann Arbor, MI 48109, dickrp@eecs.umich.edu.

This work was supported in part by NSF under grant numbers CNS-0834180, CNS07-20457, CCF-0702761, CNS-0347941, CNS-0720691, and NSF CCF-0702705, and by SRC under grant number 2007-HJ-1593. We also acknowledge the support of the U.S. Department of Education through a GAANN Fellowship for T. Chantem (award P200A090044).

run-time thermal management ever triggers throttling when there is little timing slack, the real-time task will miss its deadline. Missing hard real-time deadlines is unacceptable; an example would be failure to stop an automatically controlled train on time [22]. To guarantee hard real-time performance, designers should consider thermal effects by explicitly optimizing peak temperature while meeting all functionality and real-time deadlines. This motivates our work on the temperature-aware real-time MPSoC assignment and scheduling problem.

Existing power-aware techniques, such as energy minimization, peak power minimization, as well as global dynamic voltage and frequency scaling (DVFS), cannot solve the temperature problem in MPSoCs because they do not consider spatial thermal variation; heat generated by an active core also affects other neighboring thermal elements, be they other cores or portions of the heatsink. The net heat flow from one thermal element to another depends on the conductance parameters and the current temperatures of these thermal elements. Ignoring spatial thermal variation can lead to unnecessarily high peak temperatures, especially for high power density chips.

A. Related Work

Researchers have only recently started work on temperature-aware high-level synthesis [25] and design space exploration [20]. The objective is usually to optimize system performance subject to a peak temperature constraint. For uniprocessor architectures, Wang and Bettati presented a reactive two-speed policy to control peak temperature [43]. To guarantee real-time deadlines, a proactive thermal management policy was later proposed [5]. Rao et al. presented an optimal processor speed control policy to maximize the work completed under a temperature constraint [34]. The thermal model was later improved by the same authors [32]. Mutapic et al. focused on energy minimization under thermal and task constraints [27]. Quan et al. presented a necessary and sufficient condition for schedulability as well as a novel scheduling algorithm for real-time applications running on processors with a temperature constraint [31]. A temperature-constrained optimization technique is valuable in maximizing application performance, but it does not necessarily increase system reliability. The lower the operating temperature, the better the system reliability.

Unfortunately, there is little research that targets peak temperature control directly. Bansal et al. were among the first to study the problem of peak temperature minimization using continuous dynamic speed scaling for uniprocessors running independent tasks [2]. Jayaseelan and Mitra presented a task sequencing technique to minimize the peak temperature for periodic real-time tasks running on a single processor [13]. Neither work considers MPSoCs nor task dependencies.

The problem of assigning and scheduling real-time tasks in systems containing multiprocessors and MPSoCs has received significant research attention. Some papers focus on meeting hard real-time constraints [11] while others aim to optimize energy consumption in the presence of timing constraints [35]. Since most real-time scheduling problem variants are \mathcal{NP} -hard, many heuristics have been proposed to solve large

problem instances with different optimality criterion [37]. Once again, the focus is usually placed on meeting the thermal constraint instead of minimizing peak temperature. For example, Rao et al. presented a method to maximize throughput by determining speeds of different cores subject to a peak temperature constraint [33]. Mulas et al. proposed a task migration algorithm that balances the loads on different cores to reduce hotspots [26]. Coskun et al. used online learning [6] and integer linear program (ILP) [7] to reduce the frequency of peak temperature constraint violations. Jung et al. used dynamic thermal management (DTM) to minimize energy while meeting a peak temperature constraint [14]. An approximation algorithm for minimizing the peak temperature of ideal processors was proposed by Chen et al. for real-time tasks with no precedence constraints [4]. In addition, a temperature-aware task assignment and voltage selection algorithm was proposed by Sun et al. for three-dimensional stacked-wafer MPSoCs [41]. However, this solution cannot be used to solve our problem since only homogeneous cores are considered and lateral thermal variation is ignored (the peak temperature of 3-D MPSoCs is strongly influenced by vertical inter-core heat flow).

Xie and Hung were the first to propose a collection of heuristics for temperature-aware processor allocation, task assignment, and scheduling [45]. However, their heuristics consider spatial or temporal thermal variation, but not both types of variations. In Section V, we show that their technique can deviate significantly from optimality.

Finally, Paci et al. claim that temperature-aware design is unnecessary in low-power embedded systems [30]. While their conclusions hold for very low-power embedded processors because on-die temperature variation is small, our results show substantial ($> 30^\circ\text{C}$ improvement) benefits from temperature-aware design for MPSoCs containing several embedded processor cores (see Section VII-A for more details), using the thermal model in Section II-B.

B. Contributions

This article makes the following main contributions. We present a mixed-integer linear programming (MILP) formulation for assigning and scheduling tasks with hard real-time constraints on an MPSoC to minimize the chip peak temperature. Our formulation considers spatial and temporal thermal variations. It relies on a *phased steady-state* thermal analysis directly integrated within the MILP formulation. This analysis produces a separate steady-state thermal profile for each power profile occurring during the schedule. Extensions for temperature-dependent leakage power modeling, DVFS, finer-grained thermal modeling, and inter-task communication modeling are given.

To solve problem instances that are large or for which the effects of heat capacitance are significant, we propose a heuristic task assignment and scheduling framework in which the actual method for computing the thermal profile can be selected as appropriate. Specifically, phased steady-state thermal analysis is used when task durations are long relative to the time constants of the cores. *Transient thermal analysis*, in which temperatures are time-dependent, is used otherwise.

To exploit slack in the system where the effects of heat capacity are significant, we use the concept of delay (i.e., idle time) insertion in our transient analysis based heuristic to further reduce the chip peak temperature while guaranteeing hard real-time deadlines.

C. Organization

The paper is organized as follows. In Section II, we introduce our system model, state our assumptions, and formally define the problem. We motivate the need for a temperature-aware assignment and scheduling algorithm in Section III. We describe our formal approach in Section IV and present our flexible heuristic framework in Section V. We introduce and incorporate the concept of delay insertion into our heuristic framework in Section VI. The benefits and efficiency of our approach are experimentally determined in Section VII. Section VIII concludes the paper.

II. SYSTEM MODEL AND PROBLEM DEFINITION

The system model and the temperature-aware real-time MP-SoC assignment and scheduling problem are now described.

A. Task Model

In our model, J represents the set of hard real-time tasks to be executed. For each task $j \in J$, the worst-case execution time when running on core m is denoted by $E(j, m)$, the deadline by $D(j)$, and the release time by $R(j)$. Note that $R(j) = 0$ and $D(j) = \infty$ if no release time and deadline constraints are associated with task j . A directed acyclic graph (DAG) is used to capture data dependencies among tasks. In a DAG, nodes represent tasks and directed edges indicate data dependencies between pairs of tasks. Let Γ_{j_1, j_2} denotes the dependency between tasks j_1 and j_2 where

$$\Gamma_{j_1, j_2} = \begin{cases} 1 & \text{if task } j_1 \text{ immediately precedes task } j_2 \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

A task j may execute only after all its predecessor tasks have completed and j has been released, i.e., the current time is greater than or equal to $R(j)$. For now, we assume that there is no cost for communication among dependent tasks. This assumption will be relaxed in Section IV-E. For periodic systems, we guarantee schedule validity by scheduling out to the hyperperiod of all tasks [17]. The *hyperperiod* is the least common multiple of the periods of all tasks in the problem specification.

B. Thermal Model

We model an MPSoC with a set of cores, M . For each core $m \in M$, its width, height, and location are specified. Based on the floorplan, the set of neighbors of core m , N_m , thermal conductance to a neighbor n , $G_n(m, n)$, and thermal conductance to the heatsink element above it, $G_h(m)$, can be calculated. For each task and core combination, $P(j, m)$ indicates the power consumption of core m when executing task j . We discuss an extension to this power model to account for leakage power in Section IV-B.

Thermal analysis estimates heat transfer through heterogeneous materials among heat producers (e.g., transistors) and heat consumers (e.g., heatsinks attached to an MPSoC). In the task assignment and scheduling phase, we will adopt a coarse-grained discrete heat flow model analogous to widely used compact models [38] to balance thermal analysis efficiency and accuracy. However, the algorithm framework proposed in Section V can be used with any thermal analysis technique.

In our thermal model, which is based on the classical Fourier heat flow model, each core corresponds to a discrete thermal element; Section IV-D discusses how our approach can be modified to support finer-grained thermal element modeling. The heatsink on top of the cores is modeled using multiple thermal elements and its partitioning corresponds to the layout of the cores. Since the heatsink is usually larger than the processor itself, we model heatsink overhang using additional thermal elements; the heatsink overhangs the chip by 25% of its length and width. The interface layer is included within the heatsink instead of being modeled explicitly. The interface material is usually very thin so lateral heat flow within it can be neglected. Lateral heat flow between cores and heatsink elements is modeled.

To perform thermal analysis, we take advantage of the well-known duality between electrical and thermal circuits. The temperature of each thermal element can be expressed as a function of its power consumption, the ambient temperature, and the temperatures of neighboring thermal elements. Figure 1(a) depicts the circuit representation of this model. Here, T_A denotes the ambient temperature, $G_A(h)$ is the conductance from the heatsink element h to the ambient, and $G_{nh}(h, g)$ is the conductance between heatsink elements h and g . The current source P_m denotes the power consumption of core m . The terms $G_h(m)$ and $G_n(m, n)$ were as defined previously.

The temperature of core m at time t , $T(t, m)$, can be determined using the node thermal analysis of the circuit in Figure 1(a):

$$0 = \sum_{n \in N_m} (T(t, m) - T(t, n)) \cdot G_n(m, n) + C(m) \cdot \frac{dT(t, m)}{dt} + (T(t, m) - T(t, h)) \cdot G_h(m) - \sum_{j \in J} \alpha(t, j, m) \cdot P(j, m), \quad (2)$$

$$0 = \sum_{g \in N_h} (T(t, h) - T(t, g)) \cdot G_{nh}(h, g) + C(h) \cdot \frac{dT(t, h)}{dt} + (T(t, h) - T(t, m)) \cdot G_h(m) + (T(t, h) - T_A) \cdot G_A(h), \quad (3)$$

where $T(t, h)$ is the temperature of the heatsink element h directly above core m at time t , $C(m)$ is the thermal capacitance of core m , $C(h)$ is the thermal capacitance of heatsink element h , and $\alpha(t, j, m) = 1$ if task j is active on core m at time t . In the above two equations, if the heatsink element h is a heatsink overhang element, then the term $(T_h - T_m) \cdot G_H(m)$ is set to 0.

The thermal conductance of core m to the heatsink element h directly above it, $G_h(m)$, can be computed as described by

Serway [36]:

$$G_h(m) = \frac{Area_m}{R_{chip} \cdot Area_{chip}}, \quad (4)$$

where $Area_m$ denotes the area of core m , $Area_{chip}$ represents the area of the chip, and $R_{chip} = \frac{th_{si}}{K_{si} \cdot Area_{chip}}$, th_{si} is the thickness of silicon, and K_{si} denotes its thermal conductivity. In our experiments, we set th_{si} and K_{si} to be 0.6 mm and 148 W/mK, respectively.

The conductance of a heatsink element h to the ambient can be calculated in a similar manner. That is, we substitute $Area_m$ and $Area_{chip}$ in Eq. 4 by the area of the heatsink element under consideration and the area of the entire heatsink, respectively. In addition, we replace R_{chip} with R_{HS} in Eq. 4, where $R_{HS} = \frac{T_{active} - T_{ambient}}{P_{chip}} - R_{chip}$, P_{chip} being the total power consumption of the chip and T_{active} the average active layer temperature when all cores are busy and $T_{ambient}$ being the ambient temperature. We set T_{active} and $T_{ambient}$ to 90 °C and 45 °C, respectively.

We compute the conductance between core m and its neighbor core n as follows:

$$G_n(m, n) = \frac{w_{mn} \cdot th_{si} \cdot K_{si}}{L_{mn}}, \quad (5)$$

where w_{mn} is the length of intersection between cores m and n and L_{mn} is the distance between the midpoint of m and that of n . The lateral conductance between two heatsink elements can be computed in a similar fashion. We assume that the heatsink is made of copper, with a thickness of 1 mm and thermal conductivity of 400 W/mK.

C. Problem Definition

Given the floorplan of a chip containing a set of cores, M , and a set of hard real-time tasks, J , as described above, determine a static assignment of tasks to cores and a static, non-preemptive schedule of tasks on the cores such that all precedence constraints and real-time deadlines are met and the chip peak temperature, T_{max} , is minimized.

III. MOTIVATING EXAMPLE

Since average power (i.e., energy) for a fixed duration and peak power are related to chip temperature, it is natural to question whether optimizing peak temperature can produce significantly different results than optimizing peak power or average power. Let us consider a task set containing two identical tasks, j_1 and j_2 , each with a deadline of 5 ms. For this example, the MPSoC is arranged as shown in Figure 1(b) (core sizes are not necessarily drawn to scale in the diagram). Task execution times (E) and associated power consumptions (P) are shown near the respective cores. To minimize energy, tasks j_1 and j_2 are both assigned to core m_2 . The resultant chip peak temperature is 65.30 °C. If our objective were to minimize peak power, then task j_1 would be assigned to core m_2 and task j_2 to core m_1 , also resulting in a peak temperature of 65.30 °C. However, if task j_1 were executed on core m_4 and task j_2 on core m_1 , the peak temperature would be reduced to 65.16 °C, which is about 0.14 °C cooler. This difference is the first point in the plot in Figure 2.

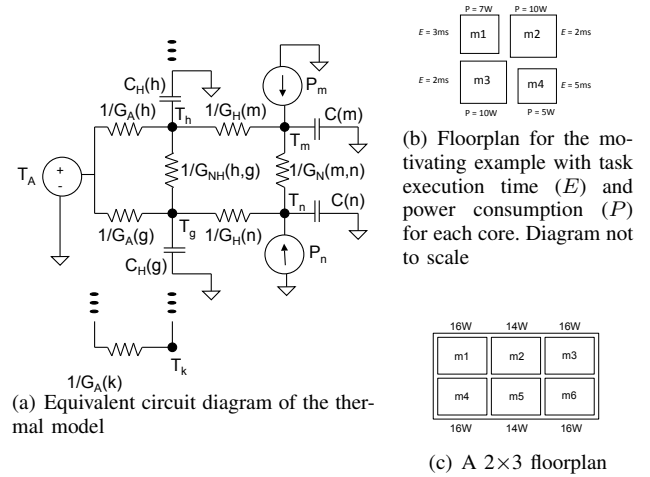


Fig. 1. Circuit and floorplans.

While the improvement in this case is small, the power density of the chip in the above example is only 0.19 W/mm². The power density can be as high as 0.79 W/mm² for 90 nm processors, 2.02 W/mm² for 65 nm processors, and 7.24 W/mm² for 45 nm processors [21]. To obtain similar chip power densities, we repeated the previous experiment but increased each core power consumption by factors of 2, 5, 10, 15, and 20. The resulting chip power densities are 0.39 W/mm², 0.97 W/mm², 1.95 W/mm², 2.92 W/mm², and 3.89 W/mm². In each case, the heatsink conductance to the ambient is adjusted to model the improved cooling solutions necessary to maintain an average active layer temperature of 90 °C.

Although task assignments and schedules are the same as before, chip peak temperatures increase when the higher power density cores are used. Figure 2 shows the reductions in chip peak temperatures when the peak temperature is optimized instead of peak power or energy for the example in Figure 1(b) and the chip power densities mentioned above. The x -axis shows the chip power density. The y -axis shows the difference between the peak temperature obtained from energy or peak power minimization and that from peak temperature minimization. As can be seen from the plot, the advantages of minimizing the chip peak temperature increase with increasing chip power density, resulting in up to 20 °C reduction in peak temperature for this example.

Energy and peak power minimization suffer from the same weakness: neither considers spatial thermal effects. In fact, energy minimization ignores both temporal and spatial thermal variation while peak power minimization considers temporal thermal variation but ignore spatial thermal variation. The peak temperature of an MPSoC is increased by crowding the same amount of energy consumption into less time and space. Hence, to minimize the chip peak temperature, tasks should be assigned and scheduled in careful consideration of thermal interaction with neighboring cores. In addition, our example indicates that although there are many cases in which minimizing peak power produces different (and potentially better) results than minimizing energy, the same results are

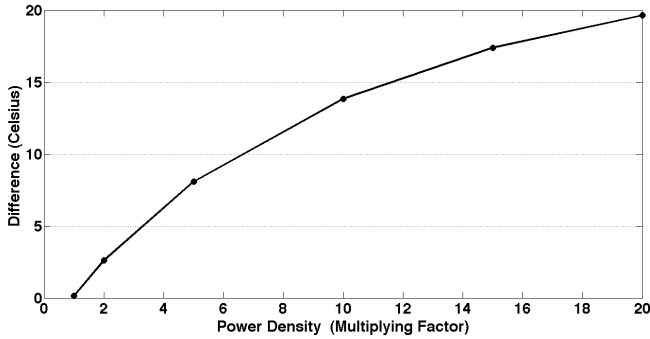


Fig. 2. Differences in peak temperatures: peak temperature minimization vs. peak power and energy minimization.

produced for some problem instances.

IV. MILP-BASED APPROACH

In this section, we present our approach to solving the problem defined in Section II-C. We also describe how our model can be extended to account for leakage power, dynamic voltage and frequency scaling (DVFS), and inter-task communication, or adjusted to use a finer-grained thermal model. Limitations of the MILP-based approach are described at the end of the section.

A. MILP Formulation

We now present our MILP formulation for the problem defined in Section II-C. We define the following variables.

$$\delta(j, m) = \begin{cases} 1 & \text{if task } j \text{ is assigned to core } m \\ 0 & \text{otherwise.} \end{cases} \quad (6)$$

$$\eta(j_1, j_2) = \begin{cases} 1 & \text{if task } j_1 \text{ starts before task } j_2 \\ 0 & \text{otherwise.} \end{cases} \quad (7)$$

$$\beta(j_1, j_2, m) = \begin{cases} 1 & \text{if task } j_2 \text{ executes on core } m, \text{ precedes}^1 \\ & \text{and overlaps with task } j_1 \\ 0 & \text{otherwise.} \end{cases} \quad (8)$$

For the sake of consistency, we let $\beta(j, j, m) \equiv \delta(j, m)$. The $\beta(j_1, j_2, m)$ variables capture the overlapping execution of different tasks and play a key role in computing the peak temperature. They are also useful in computing the peak power, as will be shown later. We use $ts(j)$ and $tf(j)$ to denote the start and finish time of task j , respectively, yielding

$$tf(j) \equiv ts(j) + \sum_{m \in M} \delta(j, m) \cdot E(j, m) \quad (9)$$

$$\equiv ts(j) + et(j). \quad (10)$$

MILP formulations have long been proposed for modeling the task assignment and scheduling problem in a heterogeneous multiprocessor environment [18]. However, energy minimization has often been the main objective. Such solutions

ignore both temporal and spatial thermal variation. Even peak power minimization only considers temporal thermal variation. To take both types of thermal variation into account, we directly minimize the chip peak temperature, T_{max} , which is the highest temperature at any position on the chip during a schedule of duration SL , i.e.,

$$T_{max} = \max_{m \in M, t \in [0, SL]} \tau_m(t), \quad (11)$$

Using Eq. 2 and Eq. 3 to compute the temperature at each node at any given time corresponds to dynamic or transient thermal analysis. Unfortunately, transient thermal analysis is computationally expensive. This makes its use in the MILP formulation impractical; the MILP solver would only be able to handle very small problem instances, thereby making it difficult to validate a heuristic. For this reason, we set the capacitance values in Eq. 2 and Eq. 3 to zero to obtain the steady-state temperature at each node when predicting temperatures in our MILP formulation. In Section IV-F, we indicate the situations in which the MILP-based approach with steady-state analysis is appropriate and inappropriate. In addition, Section V presents a solution to the more general problem of dynamic temperature optimization.

From the thermal model in Section II-B, it might appear necessary to compute the steady-state temperature, $\tau_m(t)$, of a core m at every time instant t to determine T_{max} . Even if we discretize the time duration SL , this approach may still be too costly; task execution times can vary dramatically, resulting in some tasks executing for hundreds of thousands or millions of time units. To overcome this difficulty, we make the following observations: (1) core power consumptions only change at the beginning or end of a task execution, and (2) the steady-state temperature of a core only experiences a rapid change when the power consumption of at least one core on the chip changes. Hence, we can significantly reduce the amount of computation needed to obtain T_{max} . Specifically, we only evaluate the temperature of each core m immediately after every task i starts or finishes executing on any core in the MPSoC and denote this temperature with $T(i, m)$. Consequently, the objective function of the MILP can be expressed as

$$\min T_{max}, \text{ where } T_{max} \geq T(i, m), \forall m \in M, \forall i \in J. \quad (12)$$

$T(i, m)$ satisfies the constraints given in Eq. 2 and Eq. 3, which are rewritten in Eq. 13 and Eq. 14, respectively.

$$T(i, m) \equiv T_{HS}(i, h) + \frac{1}{G_H(m)} \left[\sum_{j \in J} \beta(i, j, m) \cdot P(j, m) \right] + \frac{1}{G_H(m)} \sum_{n \in N_m} G_N(m, n) \cdot [T(i, n) - T(i, m)] \quad (13)$$

$$0 = (T_{HS}(i, h) - T(i, m)) \cdot G_H(m) + (T_{HS}(i, h) - T_A) \cdot G_A(h) + \sum_{g \in N_h} (T_{HS}(i, h) - T_{HS}(i, g)) \cdot G_{NH}(h, g). \quad (14)$$

Note that Eq. 13 is only linear if we can treat $P(j, m)$ as a constant given task j and core m . For now, we assume that

¹Precedence is not necessary but is sufficient and simplifies the test.

this is the case. We will discuss the more general case where $P(j, m)$ is not a constant in Section IV-B.

The following constraints are used to guarantee schedulability.

- 1) Every task j is assigned to exactly one core m :

$$\forall j \in J \quad \sum_{m \in M} \delta(j, m) = 1. \quad (15)$$

- 2) Every task j meets its deadline:

$$\forall j \in J \quad ts(j) + et(j) \leq d(j). \quad (16)$$

- 3) Precedence constraints are honored:

$$\forall i, j \in J \quad ts(j) \geq tf(i) \cdot \Gamma_{i,j}. \quad (17)$$

- 4) All tasks execute for their durations without overlap:

$$\forall j_1, j_2 \in J \quad 1 \leq \eta(j_1, j_2) + \eta(j_2, j_1), \quad (18)$$

$$ts(j_1) \leq ts(j_2) + (1 - \eta(j_1, j_2)) \cdot \Lambda, \quad (19)$$

$$ts(j_2) \leq ts(j_1) + \eta(j_1, j_2) \cdot \Lambda, \quad (20)$$

$$\forall j_1, j_2 \in J, j_1 \neq j_2, \forall m \in M \quad tf(j_1) \leq (2 - \delta(j_1, m) - \delta(j_2, m)) \cdot \Lambda + ts(j_2) + \Lambda \cdot (1 - \eta(j_1, j_2)), \quad (21)$$

$$tf(j_2) \leq (2 - \delta(j_1, m) - \delta(j_2, m)) \cdot \Lambda + ts(j_1) + \Lambda \cdot \eta(j_1, j_2), \quad (22)$$

where Λ is a constant greater than or equal to the largest deadline in the task set. Eq. 19 states that task j_1 must start before task j_2 if $\eta(j_1, j_2) = 1$. Eq. 21 guarantees that task j_1 finishes before task j_2 starts if tasks j_1 and j_2 are executed on the same processor and task j_1 precedes task j_2 . Similar conditions hold for Eq. 20 and 22.

Consider a situation where tasks i and j execute on cores m and n , respectively. Further, task i precedes task j and their executions overlap. At the start of task i , we only need to consider the power consumption of core m . However, at the start of task j , we must take into account the power consumptions of both cores to correctly compute the chip peak temperature. For this reason, we must ensure that $\beta(j_1, j_2, m) = 1$ only when $\delta(j_2, m) = 1$ and $ts(j_2) \leq ts(j_1) \leq tf(j_2) - \epsilon$, where ϵ is a small constant used to prevent imprecise floating point computations from making it appear as if contiguous tasks overlap in time. Therefore,

$$\forall m \in M, \forall j_1, j_2 \in J, j_1 \neq j_2 \quad tf(j_2) \geq ts(j_1) + (\beta(j_1, j_2, m) - 1) \cdot \Lambda, \quad (23)$$

$$ts(j_2) \leq ts(j_1) + (1 - \beta(j_1, j_2, m)) \cdot \Lambda, \quad (24)$$

$$1 \geq \beta(j_1, j_2, m) + \delta(j_1, m), \quad (25)$$

$$tf(j_2) - \epsilon - (1 - \eta(j_2, j_1)) \cdot \Lambda - (1 - \delta(j_2, m)) \cdot \Lambda \leq ts(j_1) + \beta(j_1, j_2, m) \cdot \Lambda. \quad (26)$$

The above MILP formulation finds an assignment and schedule that minimize the chip peak temperature. To minimize peak power P_{max} , we simply substitute the object function as follows:

$$P_{max} \geq \forall i \in J \quad \sum_{m \in M} \sum_{j \in J} \beta(i, j, m) \cdot P(j, m). \quad (27)$$

On the other hand, if total energy is to be minimized, the following objective function can be used

$$E_{total} \geq \sum_{j \in J} \sum_{m \in M} P(j, m) \cdot E(j, m) \cdot \delta(j, m), \quad (28)$$

where E_{total} denotes the total energy.

B. Modeling Power Consumption

In Eq. 13, the parameter $P(j, m)$ captures the power consumption of core m while executing task j , and

$$P(j, m) = P_{dyn}(j, m) + P_{leak}(j, m), \quad (29)$$

where P_{dyn} and P_{leak} are the dynamic power and the leakage power when running task j on core m .

Assuming average switching activity is used to evaluate $P_{dyn}(j, m)$, allows us to treat $P_{dyn}(j, m)$ as a constant. The leakage power, $P_{leak}(j, m)$, however, is a superlinear function of temperature. Simply treating $P_{leak}(j, m)$ as a constant may lead to an underestimation of the chip peak temperature. Though integrated circuit (IC) leakage power is a superlinear function of temperature, a 4-segment piecewise linear function can be used to approximate leakage in the operating temperature ranges of integrated circuits with only 0.69% error [23]. Therefore, we model the power consumption required to execute a task j on core m at temperature T_m as a piecewise linear function as follows (more segments can be added as needed)

$$P(j, m) = K_1(j, m) \cdot T_m + K_2(j, m), \quad (30)$$

where $K_1(j, m)$ and $K_2(j, m)$ are constants that depend on core m and task j . Consequently, Eq. 13 can be rewritten as

$$T(i, m) \equiv T_h(i, h) + \frac{1}{G_h(m)} \sum_{n \in N_m} G_n(m, n) \cdot [T(i, n) - T(i, m)] + \frac{1}{G_h(m)} \cdot \left[\sum_{j \in J} \beta(i, j, m) \cdot (K_1(j, m) \cdot T(i, m) + K_2(j, m)) \right]. \quad (31)$$

To eliminate the nonlinear term $\beta(i, j, m) \cdot T(i, m)$, we replace $\beta(i, j, m) \cdot T(i, m)$ with a new variable $\lambda(i, j, m)$. In other words,

$$\lambda(i, j, m) = \begin{cases} K_1(j, m) \cdot T(i, m) & \text{if } \beta(i, j, m) = 1 \\ 0 & \text{otherwise.} \end{cases} \quad (32)$$

We then add the following constraints to our MILP formulation.

$$\forall m \in M, \forall i, j \in J \quad \lambda(i, j, m) \geq 0, \quad (33)$$

$$\lambda(i, j, m) \leq \beta(i, j, m) \cdot \Lambda, \quad (34)$$

$$\lambda(i, j, m) \geq (K_1(j, m) \cdot T(i, m)) - (1 - \beta(i, j, m)) \cdot \Lambda, \quad (35)$$

$$\lambda(i, j, m) \leq (K_1(j, m) \cdot T(i, m)) - (\beta(i, j, m) - 1) \cdot \Lambda. \quad (36)$$

Solving the MILP instance given in Eq. 13, Eq. 14 (with $\beta(i, j, m) \cdot T(i, m)$ being replaced by $\lambda(i, j, m)$), Eq. 14–26, and Eq. 33–36 leads to an exact solution to the problem defined in Section II-C.

Another power consumption related issue is that we model neither inter-core interconnect nor cache power consumption. This omission is due to the following two observations. According to Ku et al. [15], the power density of the cache, and hence temperature, is relatively low due to its size. A similar conclusion can be drawn about interconnect temperatures; a simulation with default parameters using the Orion interconnect network simulator [28] revealed that typical core power density [21] is about $6.5\times$ that of a high-performance router and $10.6\times$ that of global interconnect wire.

C. Incorporating Dynamic Voltage Scaling

Many modern processors support dynamic voltage and frequency scaling (DVFS). Although using DVFS to minimize energy will generally also reduce peak temperature, energy minimization alone is not equivalent to peak temperature minimization. Energy minimization does not consider temporal or spatial thermal variation. It is, however, possible and beneficial to consider DVFS in conjunction with our peak temperature optimization technique. Our MILP formulation from Section IV-A can be modified as follows.

For each core m , the set of discrete voltage levels, K_m , must be specified. We redefine $E(j, k, m)$ to be the execution time of task j on core m at voltage level k and $P(j, k, m)$ to be the power consumption required to execute task j on core m at voltage level k . The binary variables $\delta(j, k, m)$ are also redefined to be 1 if task j is assigned to core m at voltage level k . Consequently, from Eq. 13,

$$\sum_{j \in J} \beta(i, j, m) \cdot P(j, m) = \sum_{k \in K} \sum_{j \in J} \nu(i, j, k, m) \cdot P(j, k, m),$$

where

$$\nu(j_1, j_2, k, m) = \begin{cases} 1 & \text{if } \delta(j_2, k, m) = 1, \beta(j_1, j_2, m) = 1 \\ 0 & \text{otherwise.} \end{cases}$$

The constraints in Eq. 13–26 can be readily modified for use in the new formulation.

D. Finer-Grained Thermal Model

The thermal model described in Section II-B can further be refined by using multiple thermal elements for each core, where each thermal element may have different power consumption and/or correspond to a particular functional unit of the core. Specifically, we redefine the variables $\delta(j, m, x)$ and $\beta(j_1, j_2, m, x)$ as follows.

$$\delta(j, m, x) = \begin{cases} 1 & \text{if task } j \text{ executes on functional unit} \\ & x \text{ of core } m \\ 0 & \text{otherwise.} \end{cases}$$

$$\beta(j_1, j_2, m, x) = \begin{cases} 1 & \text{if task } j_2 \text{ executes on functional unit} \\ & x \text{ of core } m, \text{ precedes, and overlaps} \\ & \text{with task } j_1 \\ 0 & \text{otherwise.} \end{cases}$$

The constraints in Eq. 13–26 can be readily modified for use in the new formulation.

Core power consumption may vary depending on the individual instructions being executed. However, the relative change in temperature at the functional-unit level is relatively slow due to the large thermal time constants of the functional units (i.e., the constant influencing the rate of temperature change in response to power consumption change). For this reason, the finer-grained thermal model presented here remains accurate; it is not necessary to model instruction-by-instruction variation in power consumption [1].

E. Modeling Inter-Task Communication

In some situations, communication cost for a task to send data to its successors is significant. Given that the time to send data from task i to task j using shared memory is expressed by parameter $C(i, j)$, our MILP formulation from Section IV-A can be modified to capture inter-task communication by simply substituting Eq. 17 with the following expression:

$$\forall i, j \in J \quad ts(j) \geq (tf(i) + C(i, j)) \cdot \Gamma_{i,j}. \quad (37)$$

F. Limitations of MILP-Based Approach

While the solution provided by the MILP formulation in Section IV-A is optimal, there are two main limitations to the MILP-based approach: (1) the MILP formulation cannot be used to efficiently solve large problem instances, as the problem defined in Section II-C is \mathcal{NP} -hard and (2) due to the use of steady-state thermal analysis, the MILP formulation may overestimate the chip peak temperature when task execution times are short relative to the thermal time constant of the cores. That is, steady-state analysis can be used to accurately predict the temperature when task execution times are long compared to the core time constants, and transient analysis should otherwise be used to permit more accurate temperature prediction, thereby allowing more aggressive scheduling of short tasks that do not cause temperatures to converge to steady-state values during execution.

V. ASSIGNMENT AND SCHEDULING HEURISTIC FRAMEWORK

To trade off accuracy in temperature estimation for running time, we propose an assignment and scheduling heuristic framework where either steady-state or transient analysis can be used, depending on the characteristics of the tasks under consideration. Additionally, our heuristic framework can be used in conjunction with any thermal modeling tool.

Our framework uses a binary search based approach to minimize peak temperature under functionality and timing constraints. It takes as inputs upper and lower temperature bounds, as well as the maximum number of iterations, *maxIter*. It then uses the average of the upper and lower bounds on the

peak temperatures as the target peak temperature to find an assignment and schedule. If an assignment and schedule is found while staying below the target temperature, the current target temperature will be used as the upper temperature bound for the next iteration of the binary search. Otherwise, it will be used as the lower temperature bound.

We introduce the key part of our framework: `ThermalSched`, a list scheduling [8] algorithm summarized in Algorithm 1. For a given task j , the earliest start time ($EST(j)$) and latest start time ($LST(j)$) are computed. The mobility of task j can then be calculated as the difference between $LST(j)$ and $EST(j)$. A potential challenge in computing $EST(j)$ and $LST(j)$ is that the execution time of task j is unknown prior to the selection of a core. Our solution is to use the smallest execution time of task j as given by the fastest core when computing $EST(j)$ and $LST(j)$ to maximize the mobility of task j , for all $j \in J$.

The steps for task assignment and scheduling follow. Ready tasks are ordered in a non-decreasing order of mobility. A *ready task* is a task whose predecessors have finished executing. Given a ready task j , `ThermalSched` selects the fastest available core that allows the task to meet its deadline while keeping the peak temperature below the target temperature. The fastest available core is chosen to maximize the mobility of the successors of task j , thereby improving schedulability. If no core is fast enough to execute task j by its deadline, `ThermalSched` terminates. (We ignore Lines 24–25 and 35–36 in Algorithm 1, and the variables `currentDelay` and `dMaxIter` for now. Their use will be explained in the next section.) Our search-based scheduling approach permits the use of an efficient list scheduler without global knowledge of temperature variation.

Observe that Algorithm 1 does not provide any details on computing the thermal profile (Line 17). Since predicting highly accurate thermal profiles increases time complexity, we propose two techniques based on the observations made in Section IV-F. These achieve different trade-offs between accuracy and time complexity.

A. Steady-State Analysis Based Heuristic

As explained in Section IV-F, if task execution times are long compared to the thermal time constants of the cores, steady-state analysis can usually rapidly and accurately predict the resulting chip temperature.

The steady-state thermal profile can be computed by expressing Eq. 2 and 3 for all the thermal elements as a system of linear equations of the form $A \cdot \mathbf{T} + B = 0$, and of size $|E| \times |E|$, where $|E|$ is the total number of thermal elements. Since the thermal conductance matrix A is fixed once a floorplan is given, the inverse of the matrix can be pre-computed once and the temperature matrix can be updated using a constant number of multiplications in each iteration. `ThermalSched` therefore has a time complexity in $\mathcal{O}(|J|^2 \cdot |M|^3)$. The time complexity of the **steady-state thermal analysis based heuristic (SSAB)** is in $\mathcal{O}(|J|^2 \cdot |M|^3 \cdot \text{maxIter})$.

Algorithm 1 ThermalSched($G(V, E)$, targetTemp , dMaxIter)

```

1: compute  $EST(j)$ , for all tasks // earliest start time
2: compute  $LST(j)$ , for all tasks // latest start time
3: compute  $\text{avgE}$  // average execution time over all tasks and
   cores
4:  $\text{mobility}(j) \leftarrow LST(j) - EST(j)$ , for all tasks
5:  $\text{currentTime} \leftarrow 0$ 
6:  $\text{busy}(m) \leftarrow 0$ , for all cores
7: while there are unscheduled tasks do
8:    $RT \leftarrow$  ready tasks in non-decreasing order of mobility
9:   for each  $j \in RT$  do
10:     $\text{invalidCount} \leftarrow 0$ 
11:     $\text{fastestCore} \leftarrow -1$ 
12:     $\text{bestExeTime} \leftarrow \infty$ 
13:    for each  $m \in M$  do
14:       $\delta(j, m) \leftarrow 0$ 
15:       $\text{endTime} \leftarrow E(j, m) + \text{currentTime}$ 
16:      if not  $\text{busy}(m)$  and  $\text{endTime} \leq D(j)$  then
17:        compute projected thermal profile for
        [ $\text{currentTime}, \text{nextIdleTime}$ ] //  $\text{nextIdleTime}$ 
        is the next earliest time when all cores become
        idle
18:         $\text{peakTemp} \leftarrow \max_{m \in M} T(j, m)$ 
19:        if  $\text{peakTemp} \leq \text{targetTemp}$  then
20:          if  $E(j, m) < \text{bestExeTime}$  then
21:             $\text{fastestCore} \leftarrow m$ 
22:             $\text{bestExeTime} \leftarrow E(j, m)$ 
23:             $\text{currentDelay} \leftarrow 0$ 
24:          else if  $\text{currentTime} > 0$  then
25:            [ $\text{fastestCore}, \text{bestExeTime}, \text{currentDelay}$ ]  $\leftarrow$ 
            DelayInsertion( $G(V, E)$ ,  $j$ ,  $m$ ,  $\text{currentTime}$ ,
             $\text{targetTemp}$ ,  $\text{dMaxIter}$ ,  $\text{avgE}$ )
26:          else if not  $\text{busy}(m)$  then
27:             $\text{invalidCount} \leftarrow \text{invalidCount} + 1$ 
28:        if  $\text{invalidCount} = |M|$  then
29:          return INFEASIBLE
30:        else if  $\text{fastestCore} \neq -1$  then
31:           $\delta(j, \text{fastestCore}) \leftarrow 1$  // assign  $j$  to  $\text{fastestCore}$ 
32:           $\text{ts}(j) \leftarrow \text{currentTime} + \text{currentDelay}$ 
33:           $\text{tf}(j) \leftarrow \text{ts}(j) + E(j, \text{fastestCore})$ 
34:           $\text{busy}(\text{fastestCore}) \leftarrow 1$ 
35:        if  $\text{currentDelay} > 0$  then
36:          break // allow no tasks to start executing between
           $\text{currentTime}$  and  $\text{currentTime} + \text{currentDelay}$ 
37:        update  $EST(j)$ , for all unscheduled tasks
38:        update  $\text{mobility}(j)$ , for all unscheduled tasks
39:         $\text{nextSchedPoint} \leftarrow \min\{\text{tf}(j) : \text{tf}(j) > \text{currentTime} +$ 
         $\text{currentDelay}\}$ 
40:        for each  $m \in M$  do
41:          if  $m$  becomes idle at  $\text{nextSchedPoint}$  then
42:             $\text{busy}(m) \leftarrow 0$ 
43:           $\text{currentTime} \leftarrow \text{nextSchedPoint}$ 
44: return FEASIBLE

```

B. Transient Analysis Based Heuristic

If task execution times are short, it is desirable to use transient analysis to compute the projected thermal profile, as explained in Section IV-F. Essentially, any existing thermal analysis technique can be used in our task assignment and scheduling heuristic framework. To validate our **transient analysis based heuristic (TAB)**, we will use `HotSpot` [38] in our experiments. The transient analysis based heuristic has a time complexity of $\mathcal{O}(|J|^2 \cdot |M| \cdot \text{maxIter} \cdot H)$, where H denotes the running time of `HotSpot` and depends on a number of input parameters such as task execution times and number of cores in the MPSoCs.

VI. DELAY INSERTION

Algorithm 1 always tries to schedule as many tasks as possible at every scheduling point to maximize the mobility of later tasks. One possible consequence of this greedy approach to task assignment and scheduling is that the chip peak temperature may be so close to the target temperature that no future ready task can execute without violating the target temperature bound, thus requiring a higher target temperature to find a feasible task assignment and schedule. To address this weakness in our heuristic framework, we introduce the concept of delay insertion. That is, when the chip peak temperature is at or near the target peak temperature, we delay the execution of the next ready task by introducing an idle interval before the task starts to allow the chip to cool down. This improves the probability of later tasks being scheduled without exceeding the target temperature bound.

Since steady-state thermal analysis depends on fast temperature rises and falls, an idle time (or a delay inserted) between task executions has no effect on the resulting peak temperature. On the other hand, transient thermal analysis would capture the cooling effects of delay insertions. Hence, the concept of delay insertions applies to the TAB heuristic only.

To demonstrate the potential benefits of delay insertions, we use the following illustrative example. Consider a system with 5 identical real-time tasks running on the MPSoC shown in Figure 1(b) with the associated execution time and $10\times$ the power consumption for each core. Without inserting any idle times, our TAB heuristic finds a feasible assignment and schedule with a peak temperature of 51.60°C . An algorithm capable of inserting appropriate delays would reduce the peak temperature to 49.88°C .

In the above example, delay insertion only reduces the chip peak temperature by 1.72°C because there are only 5 tasks in the system with relatively short execution times. In other words, executing these tasks on the example MPSoC does not significantly raise the chip peak temperature. If our tasks require $10\times$ the original execution times, i.e., a mean of 30 ms, then delay insertions would reduce the chip peak temperature by 3.87°C , from 66.22°C to 62.35°C .

We now explain the use of delay insertions in our heuristic framework. Whenever an attempt to schedule a task on a core fails because the target temperature bound is exceeded, `DelayInsertion` is called (Lines 24–25 of Algorithm 1). If an idle time has successfully been inserted into the schedule,

our heuristic will immediately move on to the next scheduling point by setting `currentTime` to $\min\{tf(j) : tf(j) > \text{currentTime} + \text{currentDelay}\}$ and continue the assignment and scheduling process (Lines 35–36 and Line 39 of Algorithm 1). No pending tasks are allowed to run during $[\text{currentTime}, \text{currentTime} + \text{currentDelay}]$. This not only simplifies the algorithm, but is also reasonable since when an idle time has been inserted, the current chip peak temperature is at or near the target temperature; executing another task within the interval $[\text{currentTime}, \text{currentTime} + \text{delay}]$ would likely cause the target temperature to be exceeded.

Algorithm 2 `DelayInsertion($G(V, E)$, j , m , currentTime , targetTemp , $d\text{MaxIter}$, avgE)`

```

1:  $\text{upperDelay} \leftarrow \text{avgE}$ 
2:  $\text{lowerDelay} \leftarrow 0$ 
3:  $\text{delay} \leftarrow (\text{upperDelay} + \text{lowerDelay}) / 2$ 
4:  $\text{iter} \leftarrow 0$ 
5:  $\text{oldPeakT} \leftarrow 0$ 
6:  $\text{newPeakT} \leftarrow 0$ 
7: while  $\text{iter} < d\text{MaxIter}$  do
8:   compute projected thermal profile for
      $[\text{currentTime}, \text{currentTime} + \text{delay}]$  and
      $[\text{currentTime} + \text{delay}, \text{nextIdleTime}]$  //  $\text{nextIdleTime}$ 
     is the next earliest time when all cores become idle
9:    $\text{newPeakT} \leftarrow \max_{m \in M} T(j, m)$ 
10:  if  $\text{targetTemp} > \text{newPeakT}$  and  $|\text{oldPeakT} - \text{newPeakT}| < \epsilon$  then
11:    if  $E(j, m) + \text{delay} < \text{bestExeTime}$  and  $\text{currentTime} + \text{delay} + E(j, m) \leq D(j)$  then
12:       $\text{fastestCore} \leftarrow m$ 
13:       $\text{bestExeTime} \leftarrow E(j, m)$ 
14:      return  $[\text{fastestCore}, \text{bestExeTime}, \text{delay}]$ 
15:    else
16:      return FAILURE
17:    else if  $\text{targetTemp} > \text{newPeakT}$  then
18:       $\text{upperDelay} \leftarrow \text{delay}$ 
19:    else
20:       $\text{lowerDelay} \leftarrow \text{delay}$ 
21:       $\text{oldPeakT} \leftarrow \text{newPeakT}$ 
22:       $\text{delay} \leftarrow (\text{upperDelay} + \text{lowerDelay}) / 2$ 
23:       $\text{iter} \leftarrow \text{iter} + 1$ 

```

Our delay insertion algorithm is shown in Algorithm 2. To find the appropriate idle time to insert delays without sacrificing the schedulability of future tasks, we use a binary search based approach. In each iteration, Algorithm 2 attempts to schedule the current task onto the core currently under consideration such that the resulting peak temperature does not exceed the target peak temperature. If this is possible, Algorithm 2 will keep this scheduling and assignment only if the current configuration has minimized the task finish time thus far. Hence, if an assignment and schedule exists for the current task that does not require delay insertions, that assignment and schedule will likely be selected (this design choice maximizes the mobility of future tasks). Algorithm 2 halts when the maximum number of iterations $d\text{MaxIter}$ has

been reached or an appropriate idle time has been found. The appropriate idle time is found when the current assignment and schedule do not exceed the target temperature and the chip peak temperature has converged (Line 10 of Algorithm 2).

In our implementation, the search begins by setting the upper bound on the delay to the average execution time of all task and core combinations and the lower bound delay to 0. This design choice has the following justifications: (1) if the upper bound is too large, *DelayInsertion* can be slow and (2) *DelayInsertion* is most likely invoked when the chip temperature is near the target temperature. Inserting an idle time similar to the average task execution time would allow, on average, the chip to cool down enough to permit most tasks to eventually be scheduled on the current core without sacrificing the efficiency of the heuristic.

VII. EXPERIMENTAL RESULTS

This section quantifies the benefits of our proposed approach and assess the quality of our heuristic framework.

A. Experimental Setup

In our experiments, we used the Embedded System Synthesis Benchmarks Suite (E3S) [10]. E3S contains 17 processing elements (PEs). In our experiments, we used the following 11 cores: AMD K6-2 450, AMD K6-2E 400 Mhz/ACR, AMD K6-2E+ 500 Mhz/ACR, AMD K6-III+ 550 Mhz/ACR, IBM PowerPC 405GP 266 Mhz, IBM PowerPC 750CX 500 MHz, IDT32334 100 MHz, IDT79RC32V334-150, IDT79RC64575 250 MHz, Motorola MPC555 40 MHz, and TI TMS320C6203 300 MHz. (Note that we did not use all 17 cores because for each floorplan, we attempted to use cores with similar sizes.) The E3S task sets follow the organization of the EEMBC benchmarks [10]. There are five benchmarks in total: Auto (24 tasks), Consumer (12 tasks), Networking (13 tasks), Office (5 tasks), and Telecom (30 tasks). Each benchmark represents an application, as its name indicates. Each sink task, which does not have any successors, has a hard real-time deadline.

For the E3S benchmarks, we experimented with a number of floorplans with 2×2 , 2×3 , and 3×3 core arrangements. We use the coarse-grained thermal model presented in Section II-B instead of the fine-grained thermal model in Section IV-D due to the lack of realistic benchmarks for which power profile variations within cores are known. Each benchmark has different floorplans, as specific tasks are required to run on specific cores among the 11 cores mentioned above. The specific configuration of each floorplan and the corresponding core names are provided in our technical report [40]. The chips consist of heterogeneous cores. Since cores with different power consumptions tend to have different areas, the vertical and lateral thermal conductance between neighboring cores, and between cores and heatsink elements will vary and can be computed as described in Section II-B.

We also used TGFF [9], which is a pseudo-random task graph generator, in our experiment to generate 10 additional benchmarks. For each benchmark, there are up to 5 task graphs and the total number of tasks ranges from 4 to 29 tasks (this is similar to the number of tasks in the E3S benchmarks). Each

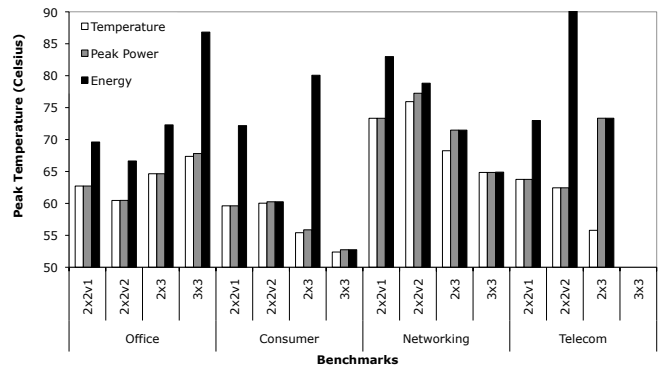


Fig. 3. Peak temperature minimization vs. energy and peak power minimization.

task has at most 3 predecessors and 2 successors. A 2×2 core arrangement was used, with an average core width and height of 5 mm and an average power consumption of 10 W.

B. MILP Formulation Performance

In this set of experiments, we used CPLEX with AMPL to solve instances of the MILP formulation in Section IV for optimal peak temperature, energy, and peak power. Each E3S benchmark was run for two 2×2 , one 2×3 , and one 3×3 floorplan.

We first examine the temperature differences between optimizing peak temperature and optimizing energy or peak power. The solutions from the MILP solver are shown in Figure 3. The x-axis shows the different benchmarks and floorplans. The y-axis shows the resulting peak temperatures. Some results are unavailable due to the MILP solver running out of memory before finding a solution. Our approach reduces peak temperatures by 9.19°C on average, and up to 24.66°C , when compared to the method that minimizes energy. Most of the improvement results from considering the effects of temporal thermal variations.

The results in Figure 3 do not show significant differences in peak temperatures between our approach and the approach that minimizes peak power. This is because the low-power embedded cores used in our benchmarks have low power densities. For example, the floorplans for the Consumer benchmarks resulted in a chip power density ranging from 0.27 W/mm^2 to 0.36 W/mm^2 with an average chip power density of 0.32 W/mm^2 . As a result, little spatial temperature variation was observed. However, spatial temperature variation will increase when higher power density chips are used, as explained in Section III.

In the first run, we used the original chip power density. The chip power density was then increased by an order of magnitude (once again to obtain similar power densities to those described by Link and Vijaykrishnan [21]) in the second run. The resulting chip power densities ranged from 0.75 W/mm^2 to 3.13 W/mm^2 with an average power density of 2.28 W/mm^2 . As shown in Figure 4, for these cores our method reduces peak temperature by 9.58°C on average, and up to 23.25°C , when compared to peak power minimization.

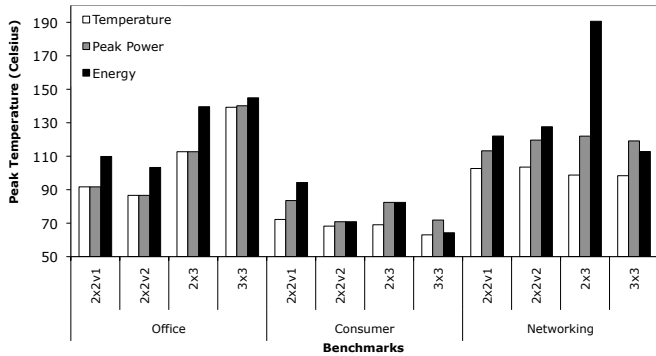


Fig. 4. Peak temperature minimization vs. energy and peak power minimization for higher power density chips.

These results demonstrate the advantage of considering spatial thermal variations.

Finally, when using the TGFF benchmarks described earlier, the MILP solver could handle eight out of ten problem instances. Minimizing peak temperature directly instead of minimizing energy reduced peak temperatures by 9.41°C on average and up to 24.19°C . In addition, when compared to peak power minimization, peak temperature minimization reduced peak temperatures by 1.27°C on average and up to 6.71°C .

The above results allow for a general conclusion to be drawn. Average power minimization suffers due to temporal and spatial thermal variation. Peak power minimization ignores spatial thermal variation. Peak temperature minimization takes both types of thermal variation into account. Variation in floorplan and other task execution times results in variation in task mobility, resulting in variation in the peak temperature improvements achieved by our approach.

C. Performance of Steady-State Analysis Based Heuristic

We assess the performance of our SSAB algorithm (Section V) by comparing its solutions to the ones from the MILP solver (Section IV-A) as well as the results from Xie's and Hung's *Heuristic 1* [45], which we refer to as the X&H heuristic. The X&H heuristic calls HotSpot to compute the temperatures. Figure 5 compares the results from the SSAB and X&H heuristics to the optimal solution from the MILP formulation. We used HotSpot to compare the peak temperatures for a fair comparison. Results for benchmarks that were not successfully solved by the X&H algorithm are omitted.

The X&H heuristic deviates from the optimal solution by 10.94°C on average and 38.40°C in the worst case. On the other hand, the SSAB heuristic finds an optimal solution in many cases while giving results that deviate by at most 3.40°C from optimality (and 0.22°C on average) requiring at most 50 binary search iterations for each benchmark. Both heuristics require similar running times, but the SSAB heuristic never performs worse than the X&H heuristic.

When the X&H heuristic was tested on the 10 TGFF benchmarks mentioned previously, we found that it could only

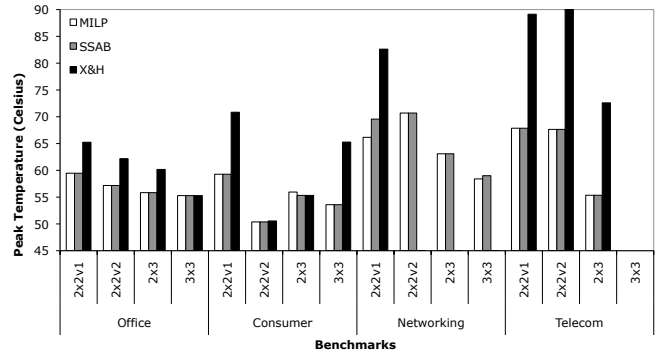


Fig. 5. Performance of steady-state analysis based heuristics (based on HotSpot).

solve two of the benchmarks, with a maximum deviation from optimality of 7.74°C . On the other hand, the SSAB heuristic could solve four, one of which was a problem instance so large that it proved intractable for the MILP solver. For the three other benchmarks, the SSAB heuristic found the same solutions as the MILP solver.

To demonstrate that the SSAB heuristic can efficiently solve larger problem instances, we considered a benchmark that consists of 30 tasks and a 4×4 core arrangement of homogeneous processors. First, we attempted to use the MILP solver. As expected, no solution was returned, as the 3.58 GB RAM workstation on which CPLEX was executing ran out of memory. The SSAB heuristic found a solution using at most 50 binary search iterations.

We used a discretized Fourier heat flow model in the SSAB heuristic. It differs from that in HotSpot. However, the experiments in this section show that the peak temperatures from the two models differed by less than 5°C on average. This indirectly served as a validation of our thermal model.

D. Performance of Transient Analysis Based Heuristic

We now assess the performance of the TAB heuristic described in Section V using the benchmarks from Section VII-A. The TAB heuristic calls HotSpot to determine transient temperatures. Since the original task execution times for the E3S benchmarks tend to be short, dynamic thermal effects can be significant. We compare the peak temperatures obtained by the MILP solver and the TAB heuristic, as shown in Figure 6. When compared to the results from the MILP solver, the TAB heuristic reduces the peak temperature by up to 0.67°C and 0.06°C on average. This is because transient analysis can more accurately predict temperatures when performing assignment and scheduling.

The TAB heuristic also improves the task finish times. Let the speedup be the ratio of the finish time of the last task in the MILP schedule to that in the TAB schedule. The maximum, minimum, and average speedups are $78.13 \times$, $1.21 \times$, and $9.02 \times$, respectively. Such a significant speedup results from the TAB heuristic being much less pessimistic in estimating temperatures and hence scheduling more tasks in parallel. However, the SSAB heuristic is more efficient

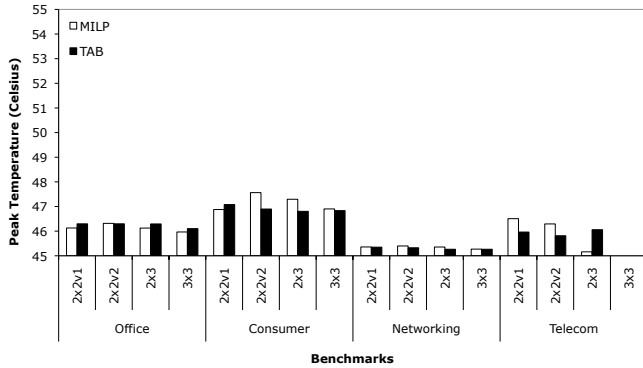


Fig. 6. Performance of transient analysis based heuristic (based on HotSpot).

than the TAB heuristic. Specifically, the SSAB heuristic is about $175\times$ faster than the TAB heuristic on average for benchmarks with short task execution times; this difference further increases for benchmarks with longer task execution times.

E. Performance of Transient Analysis Based Heuristic with Delay Insertions

To determine the impacts of delay insertions (Section VI) on reducing the chip peak temperature, we once again used the E3S and TGFF benchmarks. We compared the solutions from the original TAB heuristic to those from the improved TAB heuristic (iTAB). As before, since the power densities of the E3S cores are quite low compared to that of today's processors [21], we increased each cores power consumption by a factor of 10. The results for the E3S benchmarks are shown in Table I. The first column shows the benchmark names and associated floorplans. The second main column presents the peak temperatures from the TAB heuristic and the iTAB heuristic, as well as the differences in peak temperatures for the original task execution times. Finally, the last main column presents the data for the case where task execution times are multiplied by a factor of 10.

With the original task execution times, the effect of delay insertions is minimal for some benchmarks (e.g., networking). In fact, for two of the telecom benchmarks, iTAB actually performs worse than TAB. This is because iTAB always tries to insert delays, even when it may not be optimal to do so. For instance, if the chip is currently too hot to execute the next ready task, iTAB tries to insert idle time before scheduling that task. However, it may sometimes be better to select a different task that can meet the peak temperature constraint in Algorithm 1 without inserting delays.

When the execution times are increased by a factor of 10, we see the benefits of using the iTAB heuristic. This is because the average chip peak temperature is much higher than in the original cases and inserting idle times between task executions cools the chip down. iTAB produces solutions that reduce peak temperatures by 3.15°C on average and up to 11.92°C .

iTAB did not significantly improve on the TGFF benchmark solutions found by TAB; both solved four out of the ten

TABLE I
EFFECTIVENESS OF DELAY INSERTIONS IN REDUCING CHIP PEAK TEMPERATURES

Benchmark		Temperature ($^\circ\text{C}$)					
		Original Exe. Times (s)			Exe. Times (s) $\times 10$		
		TAB	iTAB	Diff.	TAB	iTAB	Diff.
Consumer	2 \times 2-1	65.78	60.56	5.22	84.30	80.90	3.40
	2 \times 2-2	86.24	81.47	4.77	86.24	81.47	4.77
	2 \times 3	63.32	60.06	3.26	79.24	76.81	2.43
	3 \times 3	61.97	60.28	1.69	78.60	76.85	1.75
Networking	2 \times 2-1	47.49	47.45	0.04	57.81	55.43	2.38
	2 \times 2-2	47.29	46.85	0.45	57.83	53.48	4.35
	2 \times 3	46.80	46.80	0.00	53.96	53.87	0.09
	3 \times 3	46.80	46.79	0.01	53.45	53.36	0.09
Office	2 \times 2-1	54.22	54.22	0.00	75.96	75.96	0.00
	2 \times 2-2	54.14	54.14	0.00	75.37	75.37	0.00
	2 \times 3	54.21	54.21	0.00	67.91	67.89	0.02
	3 \times 3	54.13	54.13	0.00	67.43	57.94	9.49
Telecom	2 \times 2-1	50.28	51.70	-1.42	72.06	65.53	6.53
	2 \times 2-2	49.48	52.79	-3.31	71.26	59.34	11.92
	2 \times 3	46.38	47.40	-1.02	51.37	51.30	0.08

problem instances. The largest improvement in peak temperature was 1.71°C . This is because most tasks in the TGFF benchmarks did not have enough slack for delay insertion to be effective.

Based on these results, we conclude that iTAB reduces the peak temperature of systems with time slack. It must be noted, however, that iTAB also requires longer running times. On average, iTAB takes $19.2\times$ longer to run than the original TAB algorithm. There is a trade-off between solution quality and time complexity of these algorithms.

VIII. SUMMARY AND FUTURE WORK

We presented an assignment and scheduling technique that uses a mixed-integer linear program solver to optimize IC peak temperature under precedence and hard real-time constraints based on phased steady-state thermal analysis. Experimental results show a peak temperature reduction of 10.09°C on average and up to 30.75°C for embedded processors when compared to energy minimization. When compared to peak power minimization, our approach reduced peak temperature of 8.98°C on average and up to 23.25°C for high power density chips.

To efficiently solve this \mathcal{NP} -hard assignment and scheduling problem, we also proposed a task assignment and scheduling heuristic framework in which the actual method for temperature prediction depends on task durations. Phased steady-state analysis is appropriate when task execution times are long compared to the thermal time constants of the cores and transient analysis should be used otherwise. Our phased steady-state analysis based heuristic finds an optimal solution in many cases, with a maximum deviation from optimality of 3.40°C . When compared to previous work, the heuristic achieves a temperature reduction of 10.94°C on average. The transient analysis based heuristic models and exploits the transient thermal effects of short tasks to further improve upon the existing solution by 0.67°C in the best case. Finally, we showed that incorporating the concept of delay insertion into the proposed heuristic framework results in an additional peak temperature reduction of up to 11.92°C .

Since real-time systems can experience great temperature variations at run-time due to the differences in actual task execution times, we intend to solve the peak temperature minimization problem online to further reduce temperature and increase system reliability.

ACKNOWLEDGMENT

The authors would like to thank the anonymous reviewers for their helpful suggestions.

REFERENCES

- [1] N. Allec, Z. Hassan, L. Shang, R. P. Dick, and R. Yang, "ThermalScope: multi-scale thermal analysis for nanometer-scale integrated circuits," in *Proc. Int. Conf. Computer-Aided Design*, Nov. 2008, pp. 603–610.
- [2] N. Bansal, T. Kimbrel, and K. Pruhs, "Dynamic speed scaling to manage energy and temperature," in *Proc. Symposium on Foundations of Computer Science*, Oct. 2004, pp. 520–529.
- [3] S. Y. Borkar, P. Dubey, K. C. Kahn, D. J. Kuck, H. Mulder, S. S. Pawlowski, and J. R. Rattner, "Platform 2015: Intel processor and platform evolution for the next decade," Intel Corporation, Tech. Rep., Mar. 2005.
- [4] J.-J. Chen, C.-M. Hung, and T.-W. Kuo, "On the minimization of the instantaneous temperature for periodic real-time tasks," in *Proc. Real-Time and Embedded Technology and Applications Symp.*, Apr. 2007, pp. 236–248.
- [5] J.-J. Chen, S. Wang, and L. Thiele, "Proactive speed scheduling for real-time tasks under thermal constraints," in *Proc. Real-Time and Embedded Technology and Applications Symp.*, Apr. 2009, pp. 141–150.
- [6] A. K. Coskun, T. S. Rosing, and K. Gross, "Temperature management in multiprocessor SoCs using online learning," in *Proc. Design Automation Conf.*, Jun. 2008, pp. 890–893.
- [7] A. K. Coskun, T. S. Rosing, K. Whisnant, and K. Gross, "Temperature-aware MPSoC scheduling for reducing hot spots and gradients," in *Proc. Asia & South Pacific Design Automation Conf.*, Jan. 2008, pp. 49–54.
- [8] G. De Micheli, *Synthesis and Optimization of Digital Circuits*. McGraw-Hill Book Company, NY, 1994.
- [9] R. P. Dick and D. L. R. a3snd Wayne Wolf, "TGFF: task graphs for free," in *Proc. Int. Wkshp. Hardware/Software Co-Design*, Mar. 1998, pp. 97–101.
- [10] "Embedded microprocessor benchmark consortium," <http://www.eembc.org>.
- [11] P. Gai, L. Abeni, and G. Buttazzo, "Multiprocessor DSP scheduling in system-on-a-chip architectures," in *Proc. Euromicro Conf. Real-Time Systems*, Jun. 2002, pp. 231–238.
- [12] S. H. Gunther, F. Binns, D. M. Carmean, and J. C. Hall, "Managing the impact of increasing microprocessor power consumption," *Intel Technology Journal*, vol. 5, no. 1, pp. 1–9, Feb. 2001.
- [13] R. Jayaseelan and T. Mitra, "Temperature aware task sequencing and voltage scaling," in *Proc. Int. Conf. Computer-Aided Design*, Nov. 2008, pp. 618–623.
- [14] H. Jung, P. Rong, and M. Pedram, "Stochastic modeling of a thermally-managed multi-core system," in *Proc. Design Automation Conf.*, Jun. 2008, pp. 728–733.
- [15] J. C. Ku, S. Ozdemir, G. Memik, and Y. Ismail, "Thermal management of on-chip caches through power density minimization," *IEEE Trans. VLSI Systems*, vol. 15, no. 5, pp. 592–604, May 2007.
- [16] R. Kumar, D. M. Tullsen, and N. P. Jouppi, "Core architecture optimization for heterogeneous chip multiprocessors," in *Proc. Int. Conf. Parallel Architectures and Compilation Techniques*, Sep. 2006, pp. 23–32.
- [17] E. L. Lawler and C. U. Martel, "Scheduling periodically occurring tasks on multiple processors," *Information Processing Ltrs.*, vol. 7, pp. 9–12, Feb. 1981.
- [18] L.-F. Leung, C.-Y. Tsui, and W.-H. Ki, "Simultaneous task allocation, scheduling and voltage assignment for multiple-processors-core systems using mixed integer nonlinear programming," in *Prof. Int. Symp. Circuits and Systems*, May 2003, pp. 309–312.
- [19] Y. Li, D. Brooks, Z. Hu, and K. Skadron, "Performance, energy, and thermal considerations for SMT and CMP architectures," in *Proc. Int. Symp. Computer Architecture*, Feb. 2005, pp. 71–82.
- [20] P. Lim and T. Kim, "Thermal-aware high-level synthesis based on network flow method," in *Proc. Int. Conf. Hardware/Software Codesign and System Synthesis*, Oct. 2006, pp. 124–129.
- [21] G. Link and N. Vijaykrishnan, "Thermal trends in emerging technologies," in *Proc. Int. Symp. Quality of Electronic Design*, Mar. 2006, pp. 625–632.
- [22] J. W. S. Liu, *Real-Time Systems*. Prentice-Hall, NJ, 2000.
- [23] Y. Liu, R. P. Dick, L. Shang, and H. Yang, "Accurate temperature-dependent integrated circuit leakage power estimation is easy," in *Proc. Design, Automation & Test in Europe Conf.*, Mar. 2007, pp. 1526–1531.
- [24] E. Milchman, "Intel dual-core FAQ," *Wired News*, Jul. 2006.
- [25] R. Mukherjee, S. Ögrenci Memik, and G. Memik, "Temperature-aware resource allocation and binding in high-level synthesis," in *Proc. Design Automation Conf.*, Jun. 2005, pp. 196–201.
- [26] F. Mulas, M. Pittau, M. Buttu, S. Carta, A. Acquaviva, L. Benini, and D. Atienza, "Thermal balancing policy for streaming computing on multiprocessor architectures," in *Proc. Design, Automation & Test in Europe Conf.*, Mar. 2008, pp. 734–739.
- [27] A. Mutapcic, S. Boyd, S. Murali, D. Atienza, G. D. Micheli, and R. Gupta, "Processor speed control with thermal constraints," *IEEE Trans. Circuits and Systems I*, vol. 56, no. 9, pp. 1994–2008, Sep. 2009.
- [28] "Orion 2.0: A power-performance simulator for interconnection networks," <http://www.princeton.edu/~peh/orion.html>.
- [29] "P4040: QorIQ Embedded Multicore Processor," 2009, http://www.freescale.com/webapp/sps/site/prod_summary.jsp?code=P4040&fsrc=1.
- [30] G. Paci, P. Marchal, F. Poletti, and L. Benini, "Exploring "temperature-aware design" in low-power MPSoCs," in *Proc. Design, Automation & Test in Europe Conf.*, Mar. 2006, pp. 838–843.
- [31] G. Quan, Y. Zhang, W. Wiles, and P. Pei, "Guaranteed scheduling for repetitive hard real-time tasks under the maximum temperature constraints," in *Proc. Int. Conf. Hardware/Software Codesign and System Synthesis*, Oct. 2008, pp. 267–272.
- [32] R. Rao and S. Vrudhula, "Performance optimal processor throttling under thermal constraints," in *Int. Conf. on Compilers, Architecture, and Synthesis for Embedded Systems*, Oct. 2007, pp. 257–266.
- [33] —, "Efficient online computation of core speeds to maximize the throughput of thermally constrained multi-core processors," in *Proc. Int. Conf. Computer-Aided Design*, Nov. 2008, pp. 537–542.
- [34] R. Rao, S. Vrudhula, C. Chakrabarti, and N. Chang, "An optimal analytical solution for processor speed control with thermal constraints," in *Proc. Int. Symp. Low Power Electronics & Design*, Oct. 2006, pp. 292–297.
- [35] E. Seo, Y. Koo, and J. Lee, "Dynamic repartitioning of real-time schedule on a multicore processor for energy efficiency," in *Proc. Int. Conf. Embedded and Ubiquitous Computing*, Aug. 2006, pp. 69–78.
- [36] R. A. Serway, *Physics for Scientists & Engineers with Modern Physics*. Saunders College Publishing, 1990.
- [37] G. C. Sih and E. A. Lee, "A compile-time scheduling heuristic for interconnection-constrained heterogeneous processor architectures," *IEEE Trans. Parallel & Distributed Systems*, vol. 4, no. 2, pp. 175–187, Feb. 1993.
- [38] K. Skadron, M. R. Stan, W. Huang, S. Velusamy, K. Sankaranarayanan, and D. Tarjan, "Temperature-aware microarchitecture," in *Proc. Int. Symp. Computer Architecture*, Jun. 2003, pp. 2–13.
- [39] J. Srinivasan, S. V. Adve, P. Bose, and J. A. Rivers, "Exploiting structural duplication for lifetime reliability enhancement," in *Proc. Int. Symp. Computer Architecture*, Jun. 2005, pp. 520–531.
- [40] T. Chantem, X. S. Hu, and R. P. Dick, "Temperature-aware scheduling and assignment for hard real-time applications on MPSoCs," University of Notre Dame, Tech. Rep., Dec. 2009.
- [41] C. Sun, L. Shang, and R. P. Dick, "Three-dimensional multi-processor system-on-chip thermal optimization," in *Proc. Int. Conf. Hardware/Software Codesign and System Synthesis*, Oct. 2007, pp. 117–122.
- [42] R. Viswanath, V. Wakharkar, A. Watwe, and V. Lebonheur, "Thermal performance challenges from silicon to systems," *Intel Technology Journal*, vol. 4, no. 3, pp. 1–16, Aug. 2000.
- [43] S. Wang and R. Bettati, "Delay analysis in temperature-constrained hard real-time systems with general task arrivals," in *Proc. Real-Time Systems Symp.*, Dec. 2006, pp. 323–332.
- [44] "Xbox360 Xenon," 2006, http://domino.research.ibm.com/comm/research_projects.nsf/pages/multicore.Xbox360.html.
- [45] Y. Xie and W.-L. Hung, "Temperature-aware task allocation and scheduling for embedded multiprocessor systems-on-chip (MPSoC) design," *J. VLSI Signal Processing*, vol. 45, no. 3, pp. 177–189, Dec. 2006.
- [46] T. Zhou, X. Hu, and E.-M. Sha, "Probabilistic performance estimation for real-time embedded systems," in *Int. Workshop on Timing Issues in the Specification and Synthesis of Digital Systems*, Mar. 1999, pp. 83–88.



Thidapat Chantem (S'05) received her B.S. degrees from Iowa State University in 2005 and her M.S. degree from University of Notre Dame in 2008. She is a PhD candidate in the Department of Computer Science and Engineering at University of Notre Dame. Her research interests include real-time embedded and control systems, as well as low-power and thermal-aware system-level design. She is a student member of the IEEE.



Xiaobo Sharon Hu (S'85-M'89-SM'02) received her B.S. degree from Tianjin University, China, M.S. from Polytechnic Institute of New York, and Ph.D. from Purdue University, West Lafayette, Indiana. She is Professor in the department of Computer Science and Engineering at University of Notre Dame. Her research interests include real-time embedded systems, low-power system design, VLSI and nano-scaling computing. She has published more than 160 papers in the related areas. She is currently Associate Editor for ACM Transactions on Embedded

Computing. She also served as Associate Editor for IEEE Transactions on VLSI and ACM Transactions on Design Automation of Electronic Systems and special-issue guest editor for IEEE Transactions on Industrial Informatics. She has served on the Program Committee of a number of conferences such as Design Automation Conference, International Conference on Computer-Aided Design, Design, Automation and Test in Europe Conference, IEEE Real-Time Systems Symposium, etc. She received the NSF CAREER Award in 1997, and the Best Paper Award from Design Automation Conference, 2001 and IEEE Symposium on Nanoscale Architectures, 2009.



Robert P. Dick (S'95-M'02) received his Ph.D. degree from Princeton University and his B.S. degree from Clarkson University. He was a Visiting Researcher at NEC Labs America, a Visiting Professor at Tsinghua University's Department of Electronic Engineering, and an Associate Professor at Northwestern University. He is an Associate Professor of Electrical Engineering and Computer Science at the University of Michigan. Robert received an NSF CAREER award and won his department's Best Teacher of the Year award in 2004. His technology

won a Computerworld Horizon Award and his paper was selected by DATE as one of the 30 most influential in the past 10 years in 2007. He serves on the technical program committees of several embedded systems and CAD/VLSI conferences and is an Associate Editor of IEEE Transactions on VLSI Systems.