Thermal-Aware Task Scheduling for Energy Minimization in Heterogeneous Real-Time MPSoC Systems

Junlong Zhou, Tongquan Wei, Mingsong Chen, Jianming Yan, X. Sharon Hu, and Yue Ma

Abstract-With the continuous scaling of CMOS devices, the increase in power density and system integration level have not only resulted in huge energy consumption but also led to elevated chip temperature. Thus, energy efficient task scheduling with thermal consideration has become a pressing research issue in computing systems, especially for real-time embedded systems with limited cooling techniques. In this paper, we design a twostage energy-efficient temperature-aware task scheduling scheme for heterogeneous real-time multiprocessor system-on-chip (MP-SoC) systems. In the first stage, we analyze the energy optimality of assigning real-time tasks to multiple processors of a MPSoC system, and design a task assignment heuristic that minimizes the system dynamic energy consumption under the constraint of task deadlines. In the second stage, the optimality of minimizing the peak temperature of a processor is investigated, and a slack distribution heuristic is proposed to improve the temperature profile of each processor under the thermal constraint, thus the temperature-dependent system leakage energy consumption is reduced. Through the extensive efforts made in two stages, the system overall energy consumption is minimized. Experimental results have demonstrated the effectiveness of our scheme.

Index Terms—Energy-Efficient, Thermal-Aware, Task Allocation and Scheduling, Real-Time MPSoC Systems.

I. INTRODUCTION

THE advance of technology scaling enables the integration for multiple processing elements, memory hierarchies, and dedicated hardware and I/O components on a single silicon die to form a MPSoC system. A MPSoC system is naturally heterogeneous in the sense that its processing elements such as customized hardware modules, programmable microprocessors, and embedded FPGAs have distinctive functionalities and demonstrate varying computing capability [1]. Due to their powerful parallel processing capability, higher computing density and lower clock frequencies, MPSoCs have replaced uniprocessors to become the main design paradigms for current and future embedded microprocessors in various application domains [2]. The distinct features of different types of processors of a MPSoC system can be exploited to meet the stringent design requirements of emerging real-time applications. In this paper, we focus on task scheduling issues for heterogeneous real-time MPSoC systems.

The conventional research on MPSoC systems concentrates on trading off the performance with resource requirements. Recently, increasing system integration level and decreasing

M. Chen is with the Shanghai Key Laboratory of Trustworthy Computing, East China Normal University, Shanghai 200241, China.

X. S. Hu and Y. Ma are with the Department of Computer Science and Engineering, University of Notre Dame, Notre Dame, IN 46556, USA.

feature sizes of VLSI circuits have led to a striking rise in power density [3], which not only results in huge energy consumption but also leads to elevated chip temperatures. Increase in energy consumption causes serious technical, economic, and ecological problems, such that energy management has become a critical issue in computing systems, especially for battery-powered systems that operate in harsh environments [4]–[6]. High chip temperature has adverse impact on system reliability, performance, and cost. A system will fall into the predicament of functional incorrectness, low reliability and even permanent damage if operating temperature exceeds a certain threshold [7]. Industrial studies have shown that a difference in operating temperature (10-15°C) can make a $2\times$ difference in device lifespan [8]. Thus, energy and thermal management has become a significant and pressing research issue in computing systems.

Considerable research efforts have been devoted to the investigation of task allocation for energy minimization in heterogeneous MPSoC systems. The heterogeneities of MP-SoC systems are manifested by the varying core types, different operating frequencies and power consumptions, and distinctive state switching overheads of processors. In [9], the authors addressed the problem of allocating real-time tasks onto heterogeneous cores for energy minimization under timing constraints. The presented allocation heuristics are designed as approximations to a target load distribution derived analytically. Awan et al. [10] explored the energy efficient task mapping on heterogeneous multi-core platform to reduce overall energy consumption of a real-time system. The developed heuristic first assigns tasks to processors to minimize the system active energy consumption. It then trades off higher active energy consumption for increased ability to use more efficient sleep states to reduce the system static power consumption. In [11], the authors designed a hybrid task mapping algorithm for heterogeneous MPSoCs to improve system efficiency. The hybrid method aims to maximize the throughput via static task mappings under a predefined energy budget, and further improve the performance of the mappings and reduce the energy consumption by considering the dynamic behavior of applications at runtime. All the above works attempt to fully exploit the energy saving potentials of heterogeneous processors. However, the effectiveness of utilizing the heterogeneities of a MPSoC system to reduce the chip temperature is not investigated.

Considering the temperature design constraint, Yu et al. [12] leveraged the task-level adaptability and designed a thermal-aware frequency scaling-based scheduling algorithm for maximizing the execution quality-of-service of applications on heterogeneous MPSoC platforms. The presented method

J. Zhou, T. Wei, and J. Yan are with the Shanghai Key Laboratory of Multidimensional Information Processing, and the Department of Computer Science and Technology, East China Normal University, Shanghai 200241, China. T. Wei is the Corresponding author. (email: tqwei@cs.ecnu.edu.cn.)

converts the temperature threshold into timing constraints, then optimizes the total workload cycle over all processors by judicious frequency selection. Wang et al. [13] studied the problem of reducing the peak temperature of real-time applications in MPSoC systems by utilizing system heterogeneities caused by manufacturing variations. Although it is effective to reduce the peak temperature by exploiting the heterogeneities of MPSoC systems, the energy design constraint is not discussed in these works. In addition, the heterogeneities of real-time tasks are not utilized to enhance system temperature and energy profiles.

Real-time tasks are deemed to be heterogenous when they consume different power at the same operating frequency and temperature on the same processor [14]. In [15], the heterogeneities of both system architecture and real-time tasks are used to minimize the energy consumption. A relaxationbased algorithm for three types of heterogeneous platforms are designed to achieve the task partition that is closest to the optimal solution of the relaxed problems. However, temperature is not considered as a design constraint. Saha et al. [16] developed a genetic algorithm-based task allocation that minimizes the energy consumption under the constraints of temperature limit and task deadlines. Although both energy and temperature are taken into account for optimization, the heterogeneity of real-time tasks is not considered.

In this paper, we present a static two-stage energy-efficient temperature-aware task allocation and scheduling scheme for heterogeneous real-time MPSoC systems under the constraints of task deadlines and temperature limit. The first stage of the proposed approach aims to minimize the system dynamic energy consumption by assigning the subset having a larger power dissipation factor to the processor having a smaller power dissipation factor. The second stage of the proposed approach aims to minimize the system leakage energy consumption by reducing the peak temperature of processors through slack distribution. In the two stages, feasibility analysis techniques are also designed to ensure that the target system meets its timing and thermal constraints. The major contributions of this paper are summarized as follows.

- We analyze the energy optimality of assigning tasks to multiple processors of a MPSoC system. Based on this analysis, we design a task assignment heuristic that minimizes the system dynamic energy consumption.
- We prove that the peak temperature of tasks in the thermal steady state is minimal if tasks on the same processor assume a uniform steady state temperature. Using a slack distribution policy that is developed based on this observation, the temperature profiles of processors are improved, and the temperature-dependent system leakage energy consumption is hence reduced.
- We exploit the heterogeneities of both system architecture and real-time tasks to reduce the system energy consumption. We also utilize feasibility analysis techniques to ensure real-time and temperature constraints are satisfied.

The rest of the paper is organized as follows. Section II introduces the system architecture and models, Section III shows the overview of energy minimization. Section IV presents the proposed task assignment strategy for minimizing the system dynamic energy consumption and Section V describes the proposed slack distribution policy that reduces the temperature for minimizing the system leakage energy consumption. The effectiveness of the proposed approach is verified in Section VI and concluding remarks are given in Section VII.

II. SYSTEM ARCHITECTURE AND MODELS

Consider a MPSoC system \mathcal{P} consisting of M processors $\{\mathcal{P}_1, \mathcal{P}_2, \cdots, \mathcal{P}_k, \cdots, \mathcal{P}_M\}$, where every processor \mathcal{P}_k $(1 \leq k \leq M)$ operates at a given supply voltage and processing speed pair (v_k, s_k) . Dynamic voltage scaling (DVS) is not considered in this paper since it would add another dimension for optimization [10]. In addition, the benefit of using DVS to reduce temperature is partially offset by the adverse impact of DVS on system performance.

A. Task Model

We consider real-time periodic tasks to be executed on the concerned MPSoC platform. Tasks are assumed to be heterogeneous in the sense that different tasks exhibit different power consumptions on the same processor, even executing at the same operating speed and temperature. This is due to the fact that power consumptions of tasks strongly rely on circuit activities and usage patterns of different functional units [14]. Thus, the activity factor of a task, denoted by μ (ranging in (0,1]), is introduced to capture how intensively functional units are being utilized by the task [17].

The timing characteristics of a periodic real-time task is in general described by three parameters, that is, the deadline, the period, and the worst-case execution time in cycles. A real-time task must guarantee response within a specified time constraint, which is referred to as the deadline. In a periodic real-time system, each task requires repeated execution, and the time duration between the time point of one task ready to be executed and that of the next is referred to as the period. Associating each real-time task with a worst-case execution time and a period is widely accepted in the real-time system community and is commonly adopted in the literature.

Assuming that a set Γ contains N real-time periodic tasks, denoted by $\Gamma = \{\tau_1, \tau_2, \cdots, \tau_i, \cdots, \tau_N\}$, and considering the task activity factor, the characteristics of τ_i $(1 \le i \le N)$ is described by a quadruplet $\tau_i : \{D_i, p_i, c_i, \mu_i\}$, where D_i is the deadline, p_i is the period, c_i is the worst-case execution time in cycles, and μ_i is the task activity factor. The hyper-period of set Γ , denoted by H, is the least common multiple of periods $\{p_1, p_2, \cdots, p_N\}$. Let ET(i, k) be the execution time of task τ_i on processor \mathcal{P}_k at supply voltage/speed (v_k, s_k) , that is,

$$ET(i,k) = \frac{c_i}{s_k}.$$
(1)

B. Power Model

The power consumption P of a CMOS device can be modeled as the sum of dynamic power consumption P_{dyn} and leakage (or static) power consumption P_{leak} , that is,

$$P = \hbar \cdot P_{dyn} + P_{leak}.$$
 (2)

Here \hbar is employed to represent system states and indicate whether the system is currently consuming dynamic power. Specifically, $\hbar = 1$ when the processor is in the active state and $\hbar = 0$ when the processor is in the idle state.

Dynamic power consumption mainly results from charging and discharging of gates in the circuits. It is independent of the temperature, and can be formulated as a function of supply voltage V_{dd} and operating frequency f [18], that is,

$$P_{dyn} = C^{eff} V_{dd}^2 f, aga{3}$$

where C^{eff} is the effective capacitance. Since $s \propto f$, where s is the processor speed, the power consumption of task τ_i on processor \mathcal{P}_k at the supply voltage/speed (v_k, s_k) is

$$P_{dyn}(i,k) = C_k^{eff} \mu_i v_k^2 s_k, \tag{4}$$

where μ_i is the activity factor of task τ_i .

Leakage power consumption mainly results from the leakage current and is expressed as

$$P_{leak} = N_{gate} \cdot V_{dd} \cdot I_{leak},\tag{5}$$

where N_{gate} is the number of gates, V_{dd} is the supply voltage, and I_{leak} is the leakage current. I_{leak} can be captured by a nonlinear exponential equation [19] as

$$I_{leak} = I_s (\mathcal{A}T^2 e^{(\vartheta_1 V_{dd} + \vartheta_2)/T} + \mathcal{B}e^{(\vartheta_3 V_{dd} + \vartheta_4)}), \quad (6)$$

where I_s is the leakage current at a certain reference temperature and supply voltage, T is the operating temperature, and \mathcal{A} , \mathcal{B} , ϑ_1 , ϑ_2 , ϑ_3 , and ϑ_4 are empirically determined, technology-dependent constants. (6) clearly demonstrates the complex relationship between the leakage power and temperature. However, the high-order and nonlinear terms make (6) prohibitive to perform real-time feasibility analysis. As reported in [20], the leakage current changes super linearly with the temperature and using linear approximation to model the leakage-temperature dependence can significantly simplify the leakage model while maintaining an acceptable accuracy. Therefore, as in [21], we model the leakage power of processor \mathcal{P}_k at the supply voltage/speed (v_k, s_k) as

$$P_{leak}(k) = (\alpha_k + \beta_k T) \cdot v_k, \tag{7}$$

where α_k and β_k are constants depending on processor \mathcal{P}_k .

C. Thermal Model

In an MPSoC system, each processor is assumed to be a discrete thermal element, and there is a set of heat sinks on top of the processors. These heat sinks are only used for heat dissipation and generate no power. An example layout of four processors with two heat sinks is given in Fig. 1. Heat transfer among the processors and heat sinks is a complicated dynamic process depending on the physical system. This dynamic heat transfer process can be closely approximated by Fourier's Law [16], [22]–[24], where the thermal coefficients can be obtained by using the RC models [16], [20]–[25].

Let $G_{k,m}$ represent the thermal conductance between processor \mathcal{P}_k and \mathcal{P}_m in set \mathcal{P} and $G_{k,m} = G_{m,k}$ holds for any $1 \leq k \neq m \leq M$. If there is no heat transfer between processor \mathcal{P}_k and \mathcal{P}_m , then $G_{k,m} = 0$. $G_{k,k} = 0$ holds for



Fig. 1. An example layout of four processors with two heat sinks [22].

any processor in processor set \mathcal{P} , and the thermal capacitance of processor \mathcal{P}_k is C_k . Let $\Theta = \{\Theta_1, \Theta_2, \dots, \Theta_h, \dots, \Theta_{\mathcal{H}}\}$ denote the set of \mathcal{H} heat sinks on top of the processors. The vertical thermal conductance between processor \mathcal{P}_k and heat sink Θ_h is $G_{k,h}$, which depends on the interface material and the thickness. If there is no heat dissipation from processor \mathcal{P}_k to heat sink Θ_h , then $G_{k,h} = 0$. The lateral thermal conductance between heat sink Θ_h and Θ_ℓ in set Θ is $G_{h,\ell}$ and $G_{h,\ell} = G_{\ell,h}$ holds for any $1 \le h \ne \ell \le \mathcal{H}$. The thermal conductance of a heat sink that dissipates heat to the ambient is G_{amb} . The thermal capacitance of sink Θ_h in set Θ is C_h .

Let $T_k(t)$ and $T_h(t)$ be the temperature of processor \mathcal{P}_k and heat sink Θ_h at time instance t, respectively. Let T_{amb} and $P_k(t)$ be the ambient temperature of the chip and the power consumption of processor \mathcal{P}_k at time instance t, respectively. Then according to Fourier's Law, the heat transfer process can be described as below [16], [22], [23],

$$C_k \frac{\mathrm{d}T_k(t)}{\mathrm{d}t} = P_k(t) - \sum_{\Theta_h \in \Theta} G_{k,h}(T_k(t) - T_h(t)) - \sum_{\mathcal{P}_m \in \mathcal{P}} G_{k,m}(T_k(t) - T_m(t)), \qquad (8)$$

$$C_{h} \frac{\mathrm{d}T_{h}(t)}{\mathrm{d}t} = -G_{amb}(T_{h}(t) - T_{amb})$$
$$-\sum_{\mathcal{P}_{k}\in\mathcal{P}} G_{k,h}(T_{k}(t) - T_{h}(t))$$
$$-\sum_{\Theta_{\ell}\in\Theta} G_{h,\ell}(T_{h}(t) - T_{\ell}(t)), \qquad (9)$$

where $\frac{dT_k(t)}{dt}$ and $\frac{dT_h(t)}{dt}$ are derivatives of the temperature of processor \mathcal{P}_k and heat sink Θ_h , respectively. As shown in [16], [22], [23], all these thermal parameters can be derived using the RC thermal model for a given platform.

III. OVERVIEW OF ENERGY MINIMIZATION

The focus of this work is to minimize the energy consumption of the concerned MPSoC system in a schedule duration under the constraints of real-time task deadlines and temperature limit. In this section, we first show the preliminary for estimating leakage energy consumption, then present the calculation of system overall energy consumption and define the energy minimization problem. Finally, we present the framework of our solution to solve the problem.

A. Preliminary for Leakage Energy Estimation

As the focus of this work is to minimize the overall energy consumption of the concerned MPSoC system, developing a

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/TCAD.2015.2501286, IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems



Fig. 2. An example of temperature curve with small variation in each interval.

method that can rapidly and accurately estimate the system energy consumption is of the top priority. However, this is challenging since derivation of leakage energy consumption is difficult. As introduced in Section II-B, leakage power varies with temperature and temperature is changing with time. Here the temperature refers to the operating temperature of processors since heat sinks generate no power. Either using (8) to compute the temperature at every time instance is computationally expensive or using thermal modeling tool (e.g., Hotspot [26]) to obtain the temperature profiles is time consuming. Some early works such as [27]-[30] either simply assume leakage power as a constant or totally ignore it since leakage energy consumption used to be a small part of overall energy consumption. However, with the continuous scaling of integrated circuits, the proportion of leakage in overall power dissipation is ever-increasing such that these simplistic energy models can lead to large estimation errors.

Therefore, to take into account both the accuracy and computational cost of leakage energy estimation, we adopt a compromised method that divides a schedule duration into multiple small intervals with equal length. During every interval, the operating temperature is assumed to be constant such that the leakage power consumed in the interval can be readily derived. Specifically, let [0, SD] be the schedule duration and L be the length of every interval. Then the schedule duration [0, SD] can be discretized into R intervals $[0, L], [L, 2L], \dots, [(r-1)L, rL], \dots, [(R-1)L, RL]$, where $1 \leq r \leq R$ and $R = \frac{SD}{L}$. Let $T_{k,r}^{Const}$ denote the constant operating temperature of processor \mathcal{P}_k during the interval [(r-1)L, rL], then based on (7) and our assumption, the leakage power of processor \mathcal{P}_k during [(r-1)L, rL] is

$$P_{leak}(k,r) = (\alpha_k + \beta_k T_{k,r}^{Const}) \cdot v_k.$$
(10)

This method is similar to that in [14] and is motivated by an observation illustrated in Fig. 2, that is, the temperature variation is small during each interval. Obviously, as long as the length of such interval is sufficiently small, the accuracy of this method can be very high.

When we assume the operating temperature during an interval is constant, one immediate question is what temperature should be selected for the leakage energy calculation. Since leakage is becoming the dominant source of power dissipation as the semiconductor technology advances towards the deep sub-micron era, we select the peak temperature occurring in the interval as the operating temperature. Now we only need to focus on how to obtain the peak temperature of a single interval. We take the first interval [0, L] as an example. The



4

Fig. 3. An example of execution sub-intervals.

execution of tasks on a processor during the interval [0, L] can be depicted using a sequence of execution sub-intervals, where the start time and end time of the q^{th} sub-interval is denoted by $[st_q, ed_q]$. Fig. 3 shows an example of execution sub-intervals, where three tasks are arranged to execute on a processor and eleven execution sub-intervals are produced.

It has been shown in [31] that the peak temperature can be reached at the start/end times of execution sub-intervals since the temperature during each sub-interval is monotonically increasing or decreasing. Thus, the peak temperature of tasks on processor \mathcal{P}_k in the interval [0, L] can be given as

$$T_{k,1}^{peak} = \max\{T_k(t) | t = st_1, ed_1, st_2, ed_2, \cdots, L\},\$$

where $st_1 = 0$. Since the start time of a sub-interval is the end time of its previous sub-interval, that is, $st_q = ed_{q-1}$, the peak temperature in the interval [0, L] is updated to

$$T_{k,1}^{peak} = \max\{T(t)|t = st_1, st_2, \cdots, L\}.$$
 (11)

Applying this method to the following intervals, the peak temperature of R intervals are derived, and the operating temperature of these intervals are hence obtained using that $T_{k,r}^{Const} = T_{k,r}^{peak}$ for $1 \le r \le R$. Then the calculation of leakage power is updated to

$$P_{leak}(k,r) = (\alpha_k + \beta_k T_{k,r}^{peak}) \cdot v_k.$$
(12)

Using (12), we can compute the leakage energy consumption.

B. Calculation of System Overall Energy Consumption

The N real-time tasks in set Γ are assigned to M processors in set \mathcal{P} . In other words, a given task set Γ is partitioned into M subsets $\{\Gamma_1, \Gamma_2, \dots, \Gamma_k, \dots, \Gamma_M\}$, where Γ_k is the subset of tasks assigned to processor \mathcal{P}_k . The leakage power is always consumed to maintain basic circuits and can be only eliminated by turning off the system, and the dynamic power is only consumed when executing tasks. Let E_{SD}^{tot} represent the total energy consumption of M processors in a schedule duration SD, then based on (1), (4), and (12), it can be computed as

$$E_{SD}^{tot} = \sum_{k=1}^{M} \sum_{\tau_i \in \Gamma_k} C_k^{eff} \mu_i v_k^2 s_k \cdot \frac{c_i}{s_k} \cdot \frac{SD}{p_i} + \sum_{k=1}^{M} (\alpha_k v_k SD + \beta_k v_k L \sum_{r=1}^{R} T_{k,r}^{peak}) = \sum_{k=1}^{M} (C_k^{eff} v_k^2 \sum_{\tau_i \in \Gamma_k} \frac{\mu_i c_i}{p_i}) SD + \sum_{k=1}^{M} (\alpha_k v_k SD + \beta_k v_k L \sum_{r=1}^{R} T_{k,r}^{peak}), \quad (13)$$

where the first term is the dynamic energy consumption and the second term is the static energy consumption.

5

The expression $\sum_{k=1}^{M} (C_k^{eff} v_k^2 \sum_{\tau_i \in \Gamma_k} \frac{\mu_i c_i}{p_i})$ in the first term of (13) is essentially the overall dynamic power consumption. Clearly, the dynamic energy consumption is minimal if the overall dynamic power consumption $\sum_{k=1}^{M} (C_k^{eff} v_k^2 \sum_{\tau_i \in \Gamma_k} \frac{\mu_i c_i}{p_i})$, denoted by \mathcal{V}_{metric} , is minimized. The \mathcal{V}_{metric} is in fact an energy metric to estimate the dynamic energy consumption of the MPSoC system. It can be formulated into the product of vectors, that is,

$$\begin{aligned}
\mathcal{O}_{metric}(A^M, B^M) &= A^M \times B^M \\
&= A_1 b_1 + A_2 b_2 + \dots + A_M b_M, \quad (14)
\end{aligned}$$

where $A^M = [A_1, A_2, \dots, A_k, \dots, A_M]$ and $B^M = [b_1, b_2, \dots, b_k, \dots, b_M]^T$. A^M captures processor dependent parameters, where $A_k = C_k^{eff} v_k^2$ is referred to as the power dissipation factor of processor \mathcal{P}_k . B^M captures task related parameters, where $b_k = \sum_{\tau_i \in \Gamma_k} \delta_i = \sum_{\tau_i \in \Gamma_k} \frac{\mu_i c_i}{p_i}$ is referred to as the power dissipation factor of subset Γ_k , and δ_i is the power dissipation factor of task τ_i . A^M is determined since C_k^{eff} and v_k are known for a given MPSoC system, while B^M is not determined and depends on task assignment. For a given set Γ of N real-time tasks, the sum of power dissipation factors of all tasks, denoted by $Y(\Gamma)$, can be calculated as

$$Y(\Gamma) = \sum_{i=1}^{N} \delta_i = \sum_{k=1}^{M} \sum_{\tau_i \in \Gamma_k} \frac{\mu_i c_i}{p_i} = \sum_{k=1}^{M} b_k = Y_0.$$
(15)

 $Y(\Gamma)$ is constant for a given set Γ and is denoted by Y_0 .

The expression $\sum_{k=1}^{M} (\alpha_k v_k SD + \beta_k v_k L \sum_{r=1}^{R} T_{k,r}^{peak})$ in the second term of (13) is essentially the overall leakage energy consumption in the duration SD. In this expression, α_k, β_k, v_k are constants for a given processor \mathcal{P}_k , and SD, Lare parameters decided by the scheduler. Thus, the leakage energy consumption only depends on $T_{k,r}^{peak}$, which is the peak temperature of processor \mathcal{P}_k during [(r-1)L, rL]. Obviously, the overall leakage energy consumption is minimal if the peak temperature of processors in every interval are minimized.

C. Energy Minimization Problem

As analyzed above, it is clear that the system dynamic energy consumption depends on the task assignment, and the system leakage energy consumption depends on the peak temperature of intervals. Thus, both energy-efficient task assignment and temperature-aware task scheduling are helpful to minimize the system overall energy consumption. In this paper, we propose a task assignment and scheduling scheme to address the problem of minimizing the system overall energy consumption under the real-time and thermal constraints.

Real-time constraint: In a real-time system, each task should be finished before its deadline. Suppose that the execution of real-time tasks in the system is preemptable, and the task with a smaller period has a higher priority. Let RT(i, k)denote the worst case response time of task τ_i at the supply voltage/speed (v_k, s_k) , then it can be formulated as

$$RT(i,k) = ET(i,k) + \sum_{\tau_j \in \Gamma_k, p_j < p_i} \left\lceil \frac{RT(i,k)}{p_j} \right\rceil \times ET(j,k),$$
(16)

where ET(i, k) and ET(j, k) are the execution time of task τ_i and τ_j , respectively. They both can be obtained using (1). τ_j has a higher priority than τ_i for j < i, and $\left\lceil \frac{RT(i,k)}{p_j} \right\rceil$ indicates the number of instances of τ_j during time interval RT(i, k).

Thermal constraint: The temperature of the chip should be below a temperature limit (threshold) T_{max} to avoid temperature-induced failures. The value of T_{max} is in general specified based on system design requirements. Let T_{peak} denote the peak temperature at any position on the chip during the schedule duration SD, that is,

$$T_{peak} = \max\{T(t) | \forall t \in [0, SD]\}.$$
(17)

Here T(t) can be the temperature of processors and heat sinks at time instance t. The system is deemed to be in a safe mode when the T_{peak} is below the threshold temperature T_{max} .

Problem definition: Considering the above design constraints, task allocation and scheduling problem of concerned MPSoC systems is defined as the following: Given a set Γ of N periodic real-time tasks and a set \mathcal{P} of M heterogeneous processors, derive a task allocation and scheduling scheme to minimize the system overall energy consumption in a schedule duration SD while satisfying the timing and thermal constraints. In other words, the problem can be formulated as

$$\begin{array}{ll} \mbox{Minimize:} & E_{SD}^{tot} \\ \mbox{Subject to:} & RT(i,k) \leq D_i \\ & T_{peak} \leq T_{max}. \end{array}$$

D. Framework of Our Two-Stage Solution

We propose a static two-stage task allocation and scheduling scheme to solve the above problem. As shown in Fig. 4, the proposed scheme is implemented in two stages. In the first stage, for a given task set Γ , the proposed scheme partitions tasks into M subsets and assigns them to corresponding processors, in order to minimize the system dynamic energy consumption (characterized by \mathcal{U}_{metric}). In the second stage, for the subset assigned to each processor, the proposed scheme distributes available slack on the processor to local tasks for reducing the peak temperature $T_{k,r}^{peak}$ of every interval, in order to minimize the system leakage energy consumption. Through the efforts made in two stages, the system overall energy consumption is minimized.

Feasibility analysis techniques are introduced in two stages to ensure timing and thermal constraints are met. Specifically, a real-time feasibility analysis technique is adopted in task allocation to check if the task deadlines are satisfied. A temperature feasibility analysis technique is used in task scheduling to verify the thermal constraint in the schedule duration. If the peak temperature limit is violated, the tasks that violates the thermal constraint are moved to task set for re-allocation. The proposed task assignment strategy and slack distribution policy are detailed in Section IV and V, respectively.

IV. OUR TASK ASSIGNMENT STRATEGY

This section analyzes the dynamic energy optimality of assigning tasks to multiple processors, presents a proposition on optimum task assignment, and develops a task-to-processor assignment heuristic based on the proposition. This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/TCAD.2015.2501286, IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems

6



Fig. 4. The framework of the proposed two-stage solution.

A. Analysis of the Optimality of Task-to-Processor Assignment

The system dynamic energy consumption is in fact estimated by the dynamic energy metric $\mathcal{O}_{metric}(A^M, B^M)$, which can be minimized by optimally assigning N realtime tasks to M processors. Since the vector A^M = $[A_1, A_2, \cdots, A_M]$ is constant and independent of task-toprocessor partition strategies, the dynamic energy metric $\mathcal{O}_{metric}(A^M, B^M)$ is determined by the vector $B^M = [b_1, b_2, \cdots, b_M]^T$, which varies with different task-toprocessor partition strategies. Specifically, a given set Γ of N real-time tasks can be partitioned into M subsets $\{\Gamma_1, \Gamma_2, \cdots, \Gamma_M\}$, where Γ_k $(1 \leq k \leq M)$ indicates the subset of tasks assigned to processor \mathcal{P}_k . It is clear that there are M^N instances of partitioning. In other words, assigning tasks to processors is essentially a combinatorial optimization problem. The target of combinatorial optimization problem is to find the optimum solution from all feasible solutions. Let $\Upsilon = \{\gamma_1, \gamma_2, \cdots, \gamma_n\}$ be a solution space and $f(\gamma_i)$ be the value of the objective function corresponding to the solution γ_i , then the combinatorial optimization problem involves finding the optimum solution γ^* such that $f(\gamma^*) = \min f(\gamma_i)$ holds for $\forall \gamma_i \in \Upsilon$. Since the combinatorial optimization problem is known to be NP-hard [32], the concerned task assignment problem is also NP-hard, which motivates the proposed sub-optimal task-to-processor assignment heuristic.

For a given MPSoC system, the dynamic power dissipation of processor set \mathcal{P} is characterized by a vector $A^M = [A_1, A_2, \cdots, A_M]$. For the sake of easy presentation, it is assumed that $A_1 \leq A_2 \leq \cdots \leq A_M$ holds. Similarly, the optimum power dissipation of subsets assigned to individual processors can be characterized by a vector $B^M = [B_1, B_2, \cdots, B_M]^T$, indicating that the optimum task assignment solution can minimize the objective $\mathcal{O}_{metric}(A^M, B^M)$. The sum of power dissipation factors of all assigned tasks is $Y_0 = B_1 + B_2 + \cdots + B_M$, as shown in (15). This optimal task assignment solution minimizes the dynamic energy metric $\mathcal{O}_{metric}(A^M, B^M)$ by correlating the task assignment with processor power dissipation factors, as described below.

Proposition 1: If dynamic energy metric $\mathcal{O}_{metric}(A^M, B^M)$ is minimized when $A^M = [A_1, A_2, \cdots, A_M]$ $(A_1 \leq A_2 \leq \cdots \leq A_M)$, $B^M = [B_1, B_2, \cdots, B_M]^T$, and $B_1 + B_2 + \cdots + B_M = Y_0$, then the inequality $B_1 \geq B_2 \geq \cdots \geq B_M$ holds. **Proof:** The proposition states that for an optimum task assignment solution, the processor with smaller power dissipation factor ends up with the subset of its assigned tasks having a larger power dissipation factor. As given in the proposition, the dynamic energy metric $\mathcal{O}_{metric}(A^M, B^M)$ is minimized when $A^M = [A_1, A_2, \dots, A_i, \dots, A_j, \dots, A_M]$ $(A_1 \leq A_2 \leq \dots \leq A_i \leq \dots \leq A_j \leq \dots \leq A_M),$ $B^M = [B_1, B_2, \dots, B_i, \dots, B_j \dots, B_M]^T$, and $B_1 + B_2 +$ $\dots + B_i + \dots + B_j + \dots + B_M = Y_0$, then the inequality $B_1 \geq B_2 \geq \dots \geq B_i \geq \dots \geq B_j \geq \dots \geq B_M$ holds.

Let $\mathcal{V}_{metric}(A^M, B^M)'$ be the dynamic energy metric where the position of exactly two elements in B^M is exchanged. Assume that the position of B_i and B_j (i < j)is exchanged for $\mathcal{V}_{metric}(A^M, B^M)'$, then B^M becomes $[B_1, B_2, \cdots, B_{i-1}, B_j, B_{i+1}, \cdots, B_{j-1}, B_i, B_{j+1}, \cdots, B_M]^T$ in this case. According to the definition of dynamic energy metric in (14), $\mathcal{V}_{metric}(A^M, B^M) = A_1B_1 + A_2B_2 + \cdots + A_iB_i + \cdots + A_jB_j + \cdots + A_MB_M$ and $\mathcal{V}_{metric}(A^M, B^M)' = A_1B_1 + A_2B_2 + \cdots + A_iB_j + \cdots + A_jB_i + \cdots + A_MB_M$. Since $\mathcal{V}_{metric}(A^M, B^M)$ is the optimum, $\mathcal{V}_{metric}(A^M, B^M)' - \mathcal{V}_{metric}(A^M, B^M) = (A_i - A_j)(B_j - B_i) \ge 0$. It is known that $A_i \le A_j$ for i < j, then $B_i \ge B_j$ for i < j is derived.

Given the optimum task assignment solution $B^M = [B_1, B_2, \dots, B_M]^T$ that minimizes the dynamic energy metric $\mathcal{O}_{metric}(A^M, B^M)$, any feasible solution in the solution space can be obtained by exchanging elements in $B^M = [B_1, B_2, \dots, B_M]^T$ multiple times. In each iteration of the exchange, it can be deduced that $B_i \ge B_j$ holds for i < j. In other words, the dynamic energy metric $\mathcal{O}_{metric}(A^M, B^M)$ is minimized when the processor with smaller power dissipation factor ends up with the subset of its assigned tasks having a larger power dissipation factor. The proposition is proved.

B. Task-to-Processor Assignment Heuristic

As described in Section IV-A, assigning tasks to individual processors is an NP-hard problem, which necessitates a task assignment scheme that observes the proposition presented in Section IV-A. Specifically, tasks in the subset with the maximum power dissipation factor is assigned to the processor with the minimum power dissipation factor, and tasks in the subset with the next maximum power dissipation factor is assigned to the processor with the next minimum power dissipation factor. This process repeats until all subsets of tasks are assigned to individual processors. Once a task-to-processor assignment is generated, the slack available on individual processors is distributed among local tasks. The details of the task assignment heuristic are given in Algorithm 1.

Algorithm 1 essentially partitions the tasks in the given set Γ into subsets, then assigns subsets of selected tasks to individual processors in set \mathcal{P} . The algorithm aims at minimizing the system dynamic energy consumption under the timing constraint. It is motivated by the proposition presented in Section IV-A, that is, assigning the subset having a larger power dissipation factor to the processor having a smaller power dissipation factor can minimize the system dynamic energy consumption. Since the M processors in set \mathcal{P} are sorted in the non-decreasing order of processor power dissipation factors, the

Algorithm 1: Energy-efficient task-to-processor assignment under the real-time constraint

Input: task set Γ , processor set \mathcal{P} , ambient temperature T_{amb} , and temperature limit T_{max} 1 initialization: $\{\Gamma_1, \Gamma_2, \cdots, \Gamma_M\} \leftarrow \{\emptyset, \emptyset, \cdots, \emptyset\},\$ $T_{init} \leftarrow T_{amb}$, and $k \leftarrow 1$; 2 while $\Gamma \neq \varnothing$ and $k \leq M$ do calculate task power dissipation factor δ_i of every task τ_i 3 in Γ according to $\delta_i = \frac{\mu_i c_i}{p_i}$; sort $\tau_i \in \Gamma$ in the non-increasing order of δ_i ; 4 create a temporary subset Γ_{tem} ; 5 for i = 1 to $sizeof(\Gamma)$ do /* use First-Fit to 6 group tasks into subsets */ 7 $\Gamma_{tem} = \Gamma_k + \tau_i;$ if $(RTFA(\Gamma_{tem}, k) = true)$ then 8 $\Gamma_k = \Gamma_k + \tau_i;$ 9 $\Gamma = \Gamma - \tau_i;$ 10 11 assign the slack to tasks in subset Γ_k and check the thermal constraint using Algorithm 2; $k \leftarrow k+1;$ 12 13 if $\Gamma \neq \emptyset$ and k > M then exit(1); /* the tasks in set Γ cannot be 14 feasibly scheduled under the timing and thermal constraints */ 15 else if $\Gamma = \emptyset$ and k < M then power off the vacant processors in \mathcal{P} to save energy; 16 17 **return** the target schedule $\{\Gamma_1, \Gamma_2, \cdots, \Gamma_M\}$; **Procedure** RTFA(Γ_{tem}, k) flag = true;18 19 for i = 1 to $size of(\Gamma_{tem})$ do calculate the worst case response time RT(i, k) using (16); 20 if $RT(i,k) > D_i$ then 21 flag = false; 22 break; 23 return flag; 24

focus of the algorithm becomes to derive a task-to-processor assignment that partitions tasks into M subsets, arranged in the non-increasing order of subset power dissipation factors, then assigns them to corresponding processors. This can be achieved by assigning tasks with larger task power dissipation factors to processors with smaller processor power dissipation factors. In addition, task deadlines are examined to meet the real-time constraint for each task assignment.

The pseudo code of our task assignment heuristic is given in Algorithm 1. Inputs to the algorithm are task set Γ , processor set \mathcal{P} , ambient temperature T_{amb} , and temperature limit T_{max} . Line 1 of the algorithm initializes subsets $\{\Gamma_1, \Gamma_2, \dots, \Gamma_M\}$ to $\{\emptyset, \emptyset, \dots, \emptyset\}$, chip initial temperature T_{init} to ambient temperature T_{amb} , and index k to 1. Lines 2-12 iteratively implement the process of task assignment and scheduling if the task set Γ is not empty and not all processors in \mathcal{P} have been considered. In each round of iteration, the tasks in subset Γ_k assigned to processor \mathcal{P}_k are determined in lines 3-10. More specifically, lines 3-4 calculate the task power dissipation factor δ_i of every task τ_i in set Γ and sort the tasks in the decreasing order of δ_i . Line 5 creates a temporary subset Γ_{tem} . Lines 6-10 iteratively assign tasks in set Γ to processor \mathcal{P}_k and construct subset Γ_k of tasks in a first-fit manner according to the schedulability requirement. The task with larger task power dissipation factor has higher priority when assigned to the processor. The temporary subset Γ_{tem} is used to facilitate the timing feasibility analysis of assigning task τ_i to processor \mathcal{P}_k (line 7). If the assignment can satisfy the realtime constraint, the task is assigned to the processor, and both subset Γ_k and set Γ are updated (lines 8-10). The procedure then moves to the next iteration and considers the allocation of the next task in set Γ . Otherwise, the task is not assigned and the procedure directly moves to the next iteration. The slack available on processor \mathcal{P}_k is assigned to tasks in subset Γ_k under the thermal constraint using Algorithm 2 (line 11). The process continues until a feasible schedule is generated for the system. If there is no feasible schedule for the system under the constraints, the algorithm exits (lines 13-14). When the task assignment is finished, if the system still has some vacant processors, these vacant processors are powered off for energy savings (lines 15-16). The target schedule is returned in line 17. Real-time feasibility analysis (RTFA) is called in line 8 to check if the timing design constraint is satisfied. If the response time RT(i, k) of task τ_i exceeds the deadline D_i , τ_i cannot be feasibly assigned to processor \mathcal{P}_k (lines 21-23).

V. OUR SLACK DISTRIBUTION POLICY

Real-time tasks in a given set Γ are assigned to individual processors using Algorithm 1 for reducing the dynamic energy consumption. Real-time feasibility analysis is conducted for the task assignment. Slack is the time when the processor is in the idle state, which is due to that tasks may not always take the worst-case execution time to finish and can complete earlier before the deadline. Using slack distribution can reduce processor peak temperature without increasing system dynamic or leakage energy consumption since slack distribution is in fact a rearrangement of the available slack time on the processor rather than introducing additional slacks. On the contrary, the temperature-dependent leakage energy consumption can be reduced due to the improved temperature profiles achieved by slack distribution. In this section, we focus on the design of temperature-aware slack distribution policy for minimizing the system leakage energy consumption. Thermal feasibility analysis is conducted for the slack distribution.

A. Slack Assignment to Reduce Peak Temperature

From the thermal model introduced in Section II-C and the leakage energy calculation analysis given in Section III-A, it is advantageous to do slack distribution under thermal steady state. This is because even if we discretize the schedule duration SD into a large number of small intervals, transient thermal analysis may still be too costly [23]. It has been shown in [23] that steady state thermal analysis can rapidly and accurately predict the temperature when task execution times are long compared to the thermal time constant of the processors; otherwise, it may lead to overestimated peak temperature when task execution times are short relative to the processor thermal time constants. Under the thermal steady state, we prove that the temperature profiles can be improved if all tasks on the processor assume a uniform steady state

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/TCAD.2015.2501286, IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems

8

temperature. We then discuss the policy for slack distribution which can effectively reduce the processor peak temperature.

Proposition 2: Under the thermal steady state, the peak temperature of tasks on a processor is minimal if all tasks assume a uniform steady state temperature.

Proof: Suppose that $T_{std}(i,k)$ is the steady state temperature of task τ_i on processor \mathcal{P}_k , which is formulated as [31]

$$T_{std}(i,k) = P_{std}(i,k) \times R_k + T_{amb},$$
(18)

where $P_{std}(i,k)$ is the power consumption in the steady state and can be treated as a constant since temperature becomes steady. R_k is the thermal resistance of \mathcal{P}_k and T_{amb} is the ambient temperature. Both of them are known. Thus, for the subset Γ_k , we can conclude that $\sum_{i=1}^{z} T_{std}(i,k)$ is a constant, where $z = sizeof(\Gamma_k)$. In the thermal steady state, the peak temperature of tasks are no more than their steady state temperature [23] so that the peak temperature of tasks on processor \mathcal{P}_k is max $\{T_{std}(1,k), T_{std}(2,k), \cdots, T_{std}(z,k)\}$. Since $\sum_{i=1}^{z} T_{std}(i,k)$ is a constant, it is easy to see that max $\{T_{std}(1,k), T_{std}(2,k), \cdots, T_{std}(z,k)\}$ is minimal iff $T_{std}(1,k) = T_{std}(2,k) = \cdots = T_{std}(z,k)$ holds. The proposition is proved.

The discussion above states that the peak temperature of a processor in the steady state is minimal if all tasks on the processor assume a uniform steady state temperature, which motivates the proposed slack assignment heuristic that balances steady state temperatures of tasks on the processor through slack distribution. With the improved temperature profiles, the temperature-dependent leakage energy consumption of the system is then reduced, as analyzed in Section III-B.

Let sl_i^* be the optimal slack allocated to task τ_i on processor \mathcal{P}_k for temperature balance, then the average steady power consumption $\bar{P}_{std}(i,k)$ of task τ_i during its execution time and slack time is given by

$$\bar{P}_{std}(i,k) = \frac{P_{leak}^{std}(k) \times \left(\frac{c_i}{s_k} + sl_i^*\right) + P_{dyn}(i,k) \times \frac{c_i}{s_k}}{\frac{c_i}{s_k} + sl_i^*},$$

where $P_{leak}^{std}(k)$ is the leakage power in steady state, $P_{dyn}(i,k)$ is the dynamic power, and $\frac{c_i}{s_k}$ is the task execution time.

Let $T_{std,k}$ be the uniform steady state temperature of tasks on processor \mathcal{P}_k , i.e., $T_{std}(i,k) = T_{std,k}$ holds for $\forall \tau_i \in \Gamma_k$. Given the steady state temperature and power consumption of task τ_i , the optimal slack assigned to the task can be derived by substituting $T_{std}(i,k) = T_{std,k}$ and $P_{std}(i,k) = \bar{P}_{std}(i,k)$ into (18), and is written as

$$sl_i^* = \frac{P_{dyn}(i,k) \times R_k \times c_i}{(T_{std,k} - T_{amb} - P_{leak}^{std}(k) \times R_k) \times s_k} - \frac{c_i}{s_k}.$$
 (19)

In (19), $P_{leak}^{std}(k)$ is dependent upon $T_{std,k}$, and other terms are constants either dependent upon the processor or the task, which indicates that sl_i^* is determined by $T_{std,k}$. Thus, the key of solving (19) is to derive the uniform steady state temperature of tasks on processor \mathcal{P}_k .

It has been shown in [16], [22], [23] that the uniform steady state temperature of tasks on the processor is derived when the processor reaches the steady state condition $\left(\frac{dT_k(t)}{dt} = 0\right)$. Hence we can obtain the uniform steady state temperature

 $T_{std,k}$ of each processor \mathcal{P}_k by substituting the condition $\frac{\mathrm{d}T_k(t)}{\mathrm{d}t} = 0$ into (8), which are given below [16], [22], [23].

$$\begin{bmatrix} \Omega_{1,1} & \cdots & \Omega_{1,M} \\ \Omega_{2,1} & \cdots & \Omega_{2,M} \\ \vdots & \vdots & \vdots \\ \Omega_{M,1} & \cdots & \Omega_{M,M} \end{bmatrix} \begin{bmatrix} T_{std,1} \\ T_{std,2} \\ \vdots \\ T_{std,M} \end{bmatrix} = -\begin{bmatrix} \Psi_1 \\ \Psi_2 \\ \vdots \\ \Psi_M \end{bmatrix}.$$
(20)

For any $1 \leq k \neq m \leq M$, $\Omega_{k,k} = \beta_k v_k - \sum_{h=1}^{\mathcal{H}} G_{k,h} - \sum_{m=1}^{M} G_{k,m}$, $\Omega_{k,m} = G_{k,m}$, and $\Psi_k = \alpha_k v_k + C_k^{eff} v_k^2 s_k$.

The optimal slack sl_i^* given in (19) is derived under the assumption that all the slack available on processor \mathcal{P}_k is assigned to tasks in subset Γ_k . Note that the slack assigned to a task is used to cool down the processor. Similar to the scenario that assigning slack to a task to slow down the processor increases the response time of successive tasks, assigning slack to a task to cool down the processor will lead to an increase in the response time of successive tasks. Therefore, there exists a slack bound for a task beyond which the task will miss its deadline. Let $sl_{i,max}$ be the maximum amount of slack that can be assigned to task τ_i without violating the timing constraint, then the slack actually assigned to task τ_i is given as $sl_i = \min\{sl_i^*, sl_{i,max}\}$. The next subsection describes the slack assignment heuristic in details.

B. Temperature-ware Slack Assignment Heuristic

Based on the proposed slack distribution policy, the slack assignment heuristic is developed to reduce the peak temperature of processors. The details of the heuristic are given in Algorithm 2. The algorithm iteratively assigns slack to tasks in subset Γ_k , and moves the tasks that violate thermal constraint to set Γ for re-allocation. It takes as input subset Γ_k that is generated by Algorithm 1, and an arbitrarily small positive number ϵ . In each round of iteration, the optimal slack sl_i^* of task τ_i used for temperature minimization is first calculated using (19) (line 2). Then the maximum slack $sl_{i,max}$ that could be assigned to τ_i is derived using procedure SLAK($\Gamma_k, \tau_i, \epsilon$) (line 3). The slack $sl_i = \min\{sl_i^*, sl_{i,max}\}$ is assigned to τ_i , and the task execution time is hence updated (line 4). This process repeats until all tasks in Γ_k are examined. After the slack assignment is finished, a temperature feasibility analysis (TFA) procedure is used to evaluate the thermal feasibility of the resultant task schedule, and those tasks that violate the thermal constraint are sent back to set Γ and considered to be allocated to the next processors as well as unassigned tasks.

Procedure SLAK derives the maximal slack for a task in a binary search-based manner. Inputs to the procedure are task τ_i , subset Γ_k , and the arbitrarily small positive number ϵ . A search space $[sl_l, sl_h]$ is defined and initialized to $[0, D_i - RT(i, k)]$, where sl_l and sl_h are the lower and upper bound of the space, and D_i and RT(i, k) are the deadline and response time of τ_i , respectively (line 7). The search length, denoted by ρ , is set to $sl_h - sl_l$ (line 8). Lines 9-16 describe the searching process. In each round of iteration, a dummy task τ_{tem} is created and initialized to τ_i , the median sl_{tem} of search space is calculated and taken as the slack assigned to the dummy task

 τ_{tem} , and a dummy subset Γ_{tem} is created and set to $\Gamma_k + \tau_{tem}$. Procedure RTFA presented in Algorithm 1 is called to check if timing constraint of tasks in Γ_{tem} is met. The search space $[sl_l, sl_h]$ and length ρ are updated in each iteration, and the process stops when ρ is less than yet close enough to ϵ . The lower bound sl_l of the search space is returned as the maximum slack that could be assigned to task τ_i (line 17).

The chip temperature should be below a temperature limit, as described in (17) to avoid the temperature-induced failures. To check if the thermal constraint is satisfied, we need to know the temperature of processors and heat sinks. As discussed in Section V-A, transient thermal analysis is prohibitive due to its extremely expensive computation cost, and steady state thermal analysis is less costly but may result in overestimated peak temperature. Fortunately, it is safe to use steady state thermal analysis to check the thermal constraint since if the overestimated peak temperature is below the temperature limit, the actual peak temperature must be as well. Thus, obtaining the steady state temperature of processors and heat sinks becomes the focus. At the end of Section V-A, we show the derivation of processor steady state temperatures, as in (20). Similarly, we can obtain the steady state temperature $T'_{std,h}$ of each heat sink Θ_h by substituting the steady state condition $\left(\frac{dT_h(t)}{dt} = 0\right)$ into (9), which are given below [16], [22], [23].

$$\begin{bmatrix} \Omega_{1,1}' & \cdots & \Omega_{1,\mathcal{H}}' \\ \Omega_{2,1}' & \cdots & \Omega_{2,\mathcal{H}}' \\ \vdots & \vdots & \vdots \\ \Omega_{\mathcal{H},1}' & \cdots & \Omega_{\mathcal{H},\mathcal{H}}' \end{bmatrix} \begin{bmatrix} T_{std,1}' \\ T_{std,2}' \\ \vdots \\ T_{std,\mathcal{H}}' \end{bmatrix} = -\begin{bmatrix} \Psi_1' \\ \Psi_2' \\ \vdots \\ \Psi_{\mathcal{H}}' \end{bmatrix}. \quad (21)$$

For any $1 \leq h \neq \ell \leq \mathcal{H}$, $\Omega'_{h,h} = -G_{amb} - \sum_{k=1}^{M} G_{k,h} - \sum_{\ell=1}^{\mathcal{H}} G_{h,\ell}$, $\Omega'_{h,\ell} = G_{h,\ell}$, and $\Psi'_h = G_{amb}T_{amb}$. Based on the steady state thermal analysis, we can utilize the steady state temperature to verify the thermal feasibility of the resultant task schedule, as given in Procedure TFA (lines 18-23).

VI. EVALUATION

Extensive simulation experiments have been conducted to validate the proposed scheme. The proposed scheme is compared with two benchmarking algorithms Rate Monotonic First Fit (RMFF), Rate Monotonic Best Fit (RMBF) [33], and a state-of-the-art approach Hybrid Worst-fit Genetic Algorithm (HWGA) [16]. Benchmarking algorithms RMFF and RMBF [33] are taken as baseline schemes to exhibit the energy efficiency of the proposed algorithms. The two algorithms assign priority to tasks based on task periods. A task with shorter period has higher priority than a task with longer period. RMFF is a partition algorithm that assigns the task with the highest priority to the first processor that can accommodate the task, while RMBF is a partition heuristic that assigns the task with the highest priority to the processor with smallest unused capacity among those processors on which it fits [33]. HWGA integrates a worst-fit based partition heuristic with the genetic algorithm to generate a task allocation that reduces the energy consumption while satisfying all system constraints [16]. The worst-fit based partition scheme assigns the task with the highest priority to the processor with maximum Algorithm 2: Temperature-aware slack assignment for subset Γ_k under the thermal constraint

Input: subset Γ_k , an arbitrarily small positive number ϵ **1 for** i = 1 to $sizeof(\Gamma_k)$ **do**

- $\begin{array}{c|c} \mathbf{2} & \text{calculate the optimal slack } sl_i^* \text{ of task } \tau_i \text{ using (19);} \\ \mathbf{3} & \text{derive the maximum slack } sl_{i,max} \text{ for task } \tau_i \text{ using} \end{array}$
- $\begin{array}{c|c} & SLAK(\Gamma_k, \tau_i, \epsilon); \\ \text{allocate slack } sl_i = \min\{sl_i^*, sl_{i,max}\} \text{ to } \tau_i \text{ and update} \end{array}$
- execution time $ET(i,k) = ET(i,k) + sl_i$;
- 5 if $(TFA(\Gamma_k) == false)$ then
- 6 move the tasks in Γ_k that violate the thermal constraint to Γ for re-allocation using **Algorithm** 1;

Procedure SLAK($\Gamma_k, \tau_i, \epsilon$) /* SLAK is a binary search-based method */

7 $[sl_l, sl_h] = [0, D_i - RT(i, k)];$

8 $\rho = sl_h - sl_l;$

9 while $(\epsilon < \rho)$ do

10 $| sl_{tem} = (sl_l + sl_h)/2;$ 11 $\tau_{tem} = \tau_i + sl_{tem}, \Gamma_{tem} = \Gamma_k + \tau_{te}$

11 $\tau_{tem} = \tau_i + sl_{tem}, \ \Gamma_{tem} = \Gamma_k + \tau_{tem};$ 12 if (RTFA(Γ_{tem}, k) == true) then

 $13 \quad | \quad [sl_l, sl_h] = [sl_{tem}, sl_h];$

 $\begin{array}{c|c} \mathbf{13} \\ \mathbf{14} \\ \mathbf{14} \\ \mathbf{else} \end{array} = \begin{bmatrix} s\iota_l, s\iota_h \end{bmatrix} = \begin{bmatrix} sl_{te} \\ \\ \mathbf{14} \end{bmatrix}$

15
$$|sl_l, sl_h| = [sl_l, sl_{tem}];$$

16
$$\rho = sl_h - sl_l;$$

17 return sl_l ;

```
Procedure TFA(\Gamma_k)

18 calculate the steady state temperature T_{std,k} of \mathcal{P}_k using (20);

19 calculate the steady state temperature T'_{std,h} of \Theta_h using (21);

20 if T_{std,k} \leq T_{max} and T'_{std,h} \leq T_{max} then

21 | return true;

22 else

23 | return false;
```

remaining capacity. For the sake of fair comparison, the same simulation settings are adopted for the proposed method and benchmarking algorithms RMFF, RMBF, and HWGA.

 TABLE I

 PROCESSOR PARAMETERS AND CONSTANTS [21].

v (V)	s (GHz)	α	β	C^{eff}
0.85	0.8010	7.3249	0.1666	13.0
0.90	0.8291	8.6126	0.1754	12.0
0.95	0.8553	10.238	0.1846	14.0
1.00	0.8797	12.315	0.1942	15.0
1.05	0.9027	14.998	0.2043	17.0
1.10	1.0000	18.497	0.2149	16.0

A. Experimental Settings

We perform our experimental simulations based on a 2×3 MPSoC system (M = 6). Our processor model is built on 65nm technology [19], [21]. The supply voltage v, processing speed s, and constants α , β , C^{eff} of six processors are listed in TABLE I. Four real-life benchmarks (task sets) from the Embedded System Synthesis Benchmark Suite [34] are utilized to validate the proposed scheme. The benchmarks are automotive-industrial, consumer-networking, telecom, and mpeg, which consist of 16, 20, 17, 15 tasks, respectively. As its name indicates, each benchmark represents an application. The periods of tasks in applications are assumed to equal

their deadlines. The task activity factors μ are uniformly distributed in the interval [0.4, 1], which demonstrates the heterogeneous nature of tasks [17]. We use HotSpot [26] to obtain the RC thermal model for the above platform. The floorplan and HotSpot parameters are given as follows. The number of processors is 6, the area per processor is 4mm², the die thickness is 0.15mm, the heat spreader side is 20mm, and the heat sink side is 30mm. The average of thermal resistance and capacitance of processors are selected as 0.8K/W and 340J/K, respectively. The ambient temperature is set to 45°C. Simulation experiments have been carried out under varying thermal constraints ($T_{max} = 60^{\circ}$ C, 65° C, 70° C, and 75° C) to verify the effectiveness of the proposed algorithms.

B. Simulation Results

1) Evaluation of the Accuracy and Efficiency of the Proposed Energy Estimation Method: We evaluate the accuracy and efficiency of the proposed energy estimation method, which is given in (13). Specifically, we test the performance of the proposed energy estimation method when the processor runs at different supply voltages and in schedule durations with varying lengths. The proposed energy estimation method is compared with the baseline approach presented in [14] from the aspects of energy consumption and computation cost. The baseline approach splits a schedule duration into a series of small intervals and assumes the temperature (and hence the leakage power) in every interval is close to a constant. To achieve an accurate energy estimation, we let the length of every interval be small, that is, 0.01s.

 TABLE II

 Accuracy and efficiency evaluation of the proposed energy estimation method.

SD	v	E_{Pro}	E_{Bas}	Err	ACT_{Pro}	ACT_{Bas}	Spe
(s)	(V)	(J)	(J)	(%)	(s)	(s)	(X)
5	0.9	7.45	7.18	3.6	3.8	34	8.9
	1.0	13.78	13.29	3.5			
	1.1	24.01	23.25	3.2	1		
10	0.9	15.16	14.57	3.9	4.1	46	11.2
	1.0	29.74	28.84	3.0			
	1.1	52.48	50.93	2.9			
	0.9	31.22	30.16	3.4			
20	1.0	60.57	58.94	2.7	4.0	58	14.5
	1.1	99.85	96.91	2.9			

Let E_{Pro} and E_{Bas} denote the system energy consumption calculated by using the proposed method and baseline approach [14], respectively. $Err = \frac{E_{Pro} - E_{Bas}}{E_{Pro}} \times 100\%$ denotes the relative error of the proposed method when compared to the baseline approach in terms of system energy estimation. Let ACT_{Pro} and ACT_{Bas} denote the average CPU time consumed by the proposed method and baseline approach [14], respectively. $Spe = \frac{ACT_{Bas}}{ACT_{Pro}}$ denotes the speedup achieved by the proposed method when compared to the baseline approach in terms of average CPU time. The simulation results given in TABLE II clearly show that the proposed estimation method is accurate and efficient. As can be seen in the table, the system energy consumption estimated by the proposed method is close to that of the baseline approach. The maximal relative error is no more than 3.9%. On the other hand, the proposed method can reduce the computational cost and achieves up to 14.5 times of speedup in terms of average CPU time.

2) Comparison of the Energy Consumption: We compare the proposed scheme with the methods RMFF, RMBF [33], and HWGA [16] in energy efficiency. The benchmarking methods RMFF and RMBF first arrange tasks in the order of increasing task periods, then allocate tasks to individual processors using the first fit and best fit heuristics. The stateof-the-art approach HWGA allocates tasks to individual processors using the genetic-algorithm worst-fit based heuristics. In the proposed scheme, processors are arranged in the order of increasing processor power dissipation factor while tasks are organized in the order of decreasing task power dissipation factor. Tasks with large power dissipation factor are then assigned to processors with small power dissipation factor, which has been proved to be able to minimize system dynamic energy consumption in Section IV-A. The available slack on the processor is allocated to tasks for achieving a uniform steady state temperature. Through this slack distribution, the peak temperature of tasks on the processor is minimized, as proved in Section V-A, then the temperature-dependent system leakage energy savings is hence maximized.

Fig. 5 shows the average energy consumed by the system when executing four benchmarks (automotive-industrial, consumer-networking, telecom, and mpeg) under four system thermal constraints using the proposed algorithm and three benchmarking schemes HWGA [16], RMBF, and RMFF [33]. The system thermal constraint takes the values of T_{max} = 60°C, 65°C, 70°C, and 75°C. The energy consumption given in the figure is averaged over 1000 test instances. As can be seen in the figure, the proposed algorithm consumes the least energy for a given thermal constraint among the four algorithms. Specifically, the proposed algorithm achieves energy savings of up to 11.1%, 20.1% and 23.3% as compared to benchmarking methods HWGA, RMBF, and RMFF, respectively. For example, for the scenario of benchmark automotiveindustrial under the constraint $T_{max} = 70^{\circ}$ C, the energy consumption ($E_{Pro} = 311.7$ J) of the proposed scheme is 11.1% lower than that $(E_{HWGA} = 350.6 \text{J})$ of HWGA. For the scenario of benchmark mpeg under the constraint $T_{max} =$ 75°C, the energy consumption ($E_{Pro} = 312.7$ J) of the proposed algorithm is 20.1% lower than that $(E_{RMBF} = 391.6J)$ of RMBF. For the scenario of benchmark automotive-industrial under the constraint $T_{max} = 75^{\circ}$ C, the energy consumption $(E_{Pro} = 301.5 \text{J})$ of the proposed algorithm is 23.3% lower than that $(E_{RMFF} = 393.2J)$ of RMFF.

3) Comparison of the Schedule Feasibility: We also compare the proposed scheme with benchmarking algorithms HWGA [16], RMBF, and RMFF [33] from the aspects of schedule feasibility under different thermal constraints. In addition to the algorithm adopted, the schedule feasibility is affected by two factors, that is, the input benchmark (task set) assigned to processor set and the thermal constraint set by system designer. In the simulation, we adopt four benchmarks and the temperature limits T_{max} are set to 60°C, 65°C, 70°C, and 75°C. The feasibility is calculated as the ratio of the number of benchmark instances that can be feasibly scheduled to the total number of benchmark instances. The total number



Fig. 5. The average energy consumption of benchmarks under four system thermal constraints using the proposed algorithm and three benchmarking schemes.



Fig. 6. Compare the proposed algorithm with benchmarking schemes HWGA [16], RMBF, and RMFF [33] in schedule feasibility.

of benchmark instances employed in feasibility test is 1000.

The feasibility test results are given in Fig. 6. As shown in the figure, when thermal constraint is loose ($T_{max} = 75^{\circ}$ C), the feasibility of the proposed algorithm, HWGA, RMBF, and RMFF are nearly 100%. As expected, the feasibility of the four algorithms are decreasing with the increase in benchmark size for a given thermal constraint. For instance, for the scenario of $T_{max} = 65^{\circ}$ C, the feasibility rates achieved by HWGA in the case of benchmarks automotive-industrial (N = 16), consumer-networking (N = 20), telecom (N = 17), and mpeg (N = 15) are 95.2%, 92%, 94.5%, 96% respectively. This is because increasing benchmark size leads to heavier workload on processors, which may incur violation of thermal and timing constraints. It also has been demonstrated in Fig. 6 that for a given benchmark, the feasibility of the four algorithms decreases when a rigorous thermal constraint is applied. The proposed algorithm outperforms benchmarking algorithms HWGA, RMBF, and RMFF in feasibility by up to 15%. For example, for the scenario of benchmark consumernetworking under the constraint $T_{max} = 60^{\circ}$ C, the feasibility of the proposed algorithm exceeds that of HWGA, RMBF, and RMFF by 6%, 11%, and 15%, respectively. This is primarily due to that the proposed algorithm allocates tasks to individual processors with considerations of timing and thermal constraints, while RMBF and RMFF performs task allocation without considering the thermal constraint. As compared to the state-of-the-art method HWGA, the proposed algorithm utilizes a thermal-aware slack assignment heuristic to improve processor temperature profiles and exploits feasibility analysis techniques to ensure the timeliness of the system.

4) Comparison of the Time Complexity: Due to the differences in the hardware platforms, it is difficult to directly



Fig. 7. Log-log plot of the time complexity of the proposed scheme and benchmarking methods HWGA [16], RMBF, and RMFF [33].

compare running time with the three benchmarking algorithms. Therefore, we provide a time complexity analysis in this subsection. The time complexity of the proposed scheme and benchmarking methods HWGA [16], RMBF, and RMFF [33] are $O(M^2N^2)$, $O(Max^{gen} \cdot M^2NlogM)$, O(MNlogN), and O(MNlogN), respectively, where M is the number of processors in the processor set, N is the number of tasks in the task set, and Max^{gen} is the maximum number of generations for the genetic algorithm used in method HWGA.

Fig. 7 shows the log-log plot of time complexity for the proposed scheme and benchmarking methods HWGA, RMBF, and RMFF. The plot is generated based on the setting of M = 6 and $Max^{gen} = 1000$. It has been demonstrated in the figure that the time complexity of the proposed scheme is much lower when compared to that of the state-of-the-art method HWGA, and is close to that of benchmarking methods RMBF and RMFF. The reason why RMBF and RMFF have the lowest time complexity is that they donot take into account the thermal constraint and thermal control, which adversely impacts their performance in schedule feasibility, as the results in Section VI-B3. However, the proposed scheme can not only achieve a similar low time complexity, but also has a high schedule feasibility, as the results in Section VI-B3 and VI-B4.

VII. CONCLUSION

This paper proposes a task allocation scheme and a slack assignment policy for heterogeneous real-time MPSoC systems. The proposed task allocation scheme minimizes the system dynamic energy consumption by assigning tasks to individual processors in the way that the processor with a small power dissipation factor ends up with allocated tasks in a subset having a large power dissipation factor. The proposed slack assignment policy that reduces the system leakage energy consumption by improving processor temperature profiles, is motivated by the observation that the peak temperature of tasks on a processor is minimal if tasks assume a uniform steady state temperature. Feasibility analysis techniques are utilized to ensure the timing and thermal constraints can be satisfied.

The proposed energy estimation method is evaluated with aspects to accuracy and efficiency. The simulation results show that as compared to the baseline approach, the proposed method can not only accurately estimate the energy consumption within 3.9% relative error, but also achieve up to $14.5 \times$ speedup in terms of CPU time. The proposed algorithms are

compared with benchmarking schemes RMFF, RMBF, and HWGA in terms of energy efficiency and schedule feasibility. The simulation results show that the proposed algorithms consume up to 23.3% less energy and achieve up to 15% higher feasibility as compared to benchmarking schemes.

ACKNOWLEDGEMENTS

This work was in part supported by National Natural Science Foundation of China under the grant 91418203 and 61202103. This work was also partially supported by ECNU Outstanding Doctoral Dissertation Cultivation Plan of Action under the grant PY2015047.

REFERENCES

- F. Wang, C. Nicopoulos, X. Wu, X. Xie, and N. Vijaykrishnan, "Variationaware task allocation and scheduling for MPSoC," in *Proc. Int. Conf. Computer-Aided Design*, pp. 598-603, 2007.
- [2] H. Javaid, M. Shafique, J. Henkel, and S. Parameswaran, "Energy-efficient adaptive pipelined MPSoCs for multimedia applications," *IEEE Trans. Computer-Aided Design of Integrated Circuits and Systems*, vol. 33, no. 5, pp. 663-676, 2014.
- [3] S. Pagani, J. Chen, and J. Henkel, "Energy and peak power efficiency analysis for the single voltage approximation (SVA) scheme," *IEEE Trans. Computer-Aided Design of Integrated Circuits and Systems*, 2015.
- [4] A. Ejlali, B. Al-Hashimi, and P. Eles, "Low-energy standby-sparing for hard real-time systems," *IEEE Trans. Computer-Aided Design of Integrated Circuits and Systems*, vol. 31, no. 3, pp. 329-342, 2012.
- [5] G. Chen, K. Huang, and A. Knoll, "Energy optimization for real-time multiprocessor system-on-chip with optimal DVFS and DPM combination," ACM Trans. Embedded Computing Systems, vol. 13, no. 3s, 2014.
- [6] M. Shafique, L. Bauer, and J. Henkel, "Adaptive energy management for dynamically reconfigurable processors," *IEEE Trans. Computer-Aided Design of Integrated Circuits and Systems*, vol. 33, no. 1, pp. 50-63, 2014.
- [7] J. Srinivasan, S. Adve, P. Bose, and J. Rivers, "Exploiting structural duplication for lifetime reliability enhancement," in *Proc. Int. Symp. Computer Architecture*, pp. 520-531, 2005.
- [8] R. Viswanath, W. Vijay, A. Watwe, and V. Lebonheur, "Thermal performance challenges from silicon to systems," *Intel Technology Journal*, vol. 4, no. 3, pp. 1-16, 2000.
- [9] A. Colin, A. Kandhalu, and R. Rajkumari, "Energy-efficient allocation of real-time applications onto heterogeneous processors," in *Proc. Int. Conf. Embedded and Real-Time Computing Systems and Applications*, pp. 1-10, 2014.
- [10] M. Awan and S. Petters, "Energy aware partitioning of tasks onto a heterogeneous multi-core platform," in *Proc. Int. Symp. Real-Time and Embedded Technology and Applications*, pp. 205-214, 2013.
- [11] W. Quan and A. Pimentel, "A hybrid task mapping algorithm for heterogeneous MPSoCs," ACM Trans. Embedded Computing Systems, vol. 14, no. 1, 2015.
- [12] H. Yu, R. Syed, and Y. Ha, "Thermal-aware frequency scaling for adaptive workloads on heterogeneous MPSoCs," in *Proc. Int. Conf. Design, Automation and Test in Europe*, 2014.
- [13] T. Wang, M. Fan, G. Quan, and S. Ren, "Heterogeneity exploration for peak temperature reduction on multi-core platforms," in *Proc. Int. Symp. Quality Electronic Design*, pp. 107-114, 2014.
- [14] Y. Liu, R. Dick, L. Shang, and H. Yang, "Thermal vs energy optimization for DVFS-enabled processors in embedded systems," in *Proc. Int. Symp. Quality Electronic Design*, pp. 204-209, 2007.
- [15] D. Li and J. Wu, "Minimizing energy consumption for frame-based tasks on heterogeneous multiprocessor platforms," *IEEE Trans. Parallel* and Distributed Systems, vol. 26, no. 3, pp. 810-823, 2015.
- [16] S. Saha, Y. Lu, and J. Deogun, "Thermal-constrained energy-aware partitioning for heterogeneous multi-core multiprocessor real-time systems," in *Proc. Int. Conf. on Embedded and Real-Time Computing Systems and Applications*, pp. 41-50, 2012.
- [17] H. Huang, V. Chaturvedi, G. Quan, J. Fan, and M. Qiu, "Throughput maximization for periodic real-time systems under the maximal temperature constraint," ACM Trans. Embedded Computing Systems, vol. 13, no. 2s, 2014.
- [18] N. Weste and K. Eshraghian, "Principles of CMOS VLSI design: A system perspective," Addison-Wesley Publishing Company, 1992.

13

- [19] W. Liao, L. He, and K. Lepak, "Temperature and supply voltage aware performance and power modeling at microarchitecture level," *IEEE Trans. Computer-Aided Design of Integrated Circuits and Systems*, vol. 24, no. 7, pp. 1042-1053, 2005.
- [20] Y. Liu, R. Dick, L. Shang, and H. Yang, "Accurate temperaturedependent integrated circuit leakage power estimation is easy," in *Proc. Int. Conf. Design Automation and Test in Europe*, pp. 1526-1531, 2007.
- [21] G. Quan and V. Chaturvedi, "Feasibility analysis for temperature constraint hard real-time periodic tasks," *IEEE Trans. Industrial Informatics*, vol. 6, no. 3, pp. 329-339, 2010.
- [22] N. Fisher, J. Chen, S. Wang, and L. Thiele, "Thermal-aware global realtime on multicore systems", in *Proc. Int. Symp. Real-Time and Embedded Technology and Applications*, pp. 131-140, 2009.
- [23] T. Chantem, X. Hu, and R. Dick, "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.
- [24] S. Zhang and K. Chatha, "Approximation algorithm for the temperatureaware scheduling problem", in *Proc. Int. Conf. Computer-Aided Design*, pp. 281-288, 2007.
- [25] K. Skadron, M. Stan, K. Sankaranarayanan, W. Huang, S. Velusamy, and D. Tarjan, "Temperature-aware microarchitecture: Modeling and implementation," *ACM Trans. Architecture and Code Optimization*, vol. 1, no. 1, pp. 94-125, 2004.
- [26] HotSpot. University of Virginia. [Online]. Available: http://lava.cs. virginia.edu/HotSpot.
- [27] B. Zhao, H. Aydin, and D. Zhu, "On maximizing reliability of realtime embedded applications under hard energy constraint," *IEEE Trans. Industrial Informatics*, vol. 6, no.3, pp. 316-328, 2010.
- [28] D. Zhu, R. Melhem, and B. Childers, "Scheduling with dynamic voltage/speed adjustment using slack reclamation in multiprocessor realtime systems," in *Proc. Int. Symp. Real-Time Systems*, pp. 84-94, 2001.
- [29] J. Chen, H. Hsu, and T. Kuo, "Leakage-aware energy-efficient scheduling of real-time tasks in multiprocessor systems," in *Proc. Int. Symp. Real-Time and Embedded Technology and Applications*, pp. 408-417, 2006.
- [30] K. Li, "Scheduling precedence constrained tasks with reduced processor energy on multiprocessor computers", *IEEE Trans. Computers*, vol. 61, no. 12, pp. 1668-1681, 2012.
- [31] R. Jayaseelan and T. Mitra, "Temperature aware task sequencing and voltage scaling," in *Proc. Int. Conf. Computer-Aided Design*, pp. 618-623, 2008.
- [32] B. Korte, J. Vygen, B. Korye, and J. Vygen, "Combinatorial Optimization," Springer, 2002.
- [33] O. Zapata and P. Alvarez, "EDF and RM multiprocessor scheduling algorithms: survey and performance evaluation," *Seccion de Computacion Av. IPN*, 2005.
- [34] E3S. [Online]. Available: http://ziyang.eecs.umich.edu/~dickrp/e3s/. 2013.