

A Design and Analysis Framework for Thermal-Resilient Hard-Real-Time Systems

Pradeep M. Hettiarachchi, Wayne State University

Nathan Fisher, Wayne State University

Masud Ahmed, Wayne State University

Le Yi Wang, Wayne State University

Shinan Wang, Wayne State University

Weisong Shi, Wayne State University

We address the challenge of designing predictable real-time systems in an unpredictable thermal environment where environmental temperature may dynamically change (e.g., implantable medical devices). Towards this challenge, we propose a control-theoretic design methodology which permits a system designer to specify a set of hard-real-time performance modes under which the system may operate. The system automatically adjusts the real-time performance mode based on the external thermal stress. We show (via analysis, simulations, and a hardware testbed implementation) that our control-design framework is stable and control performance is equivalent to previous real-time thermal approaches, even under dynamic temperature changes. A crucial and novel advantage of our framework over previous real-time control is the ability to guarantee hard deadlines even under transitions between modes. Furthermore, our system design permits the calculation of a new metric called *thermal resiliency* which characterizes the maximum external thermal stress that any hard-real-time performance mode can withstand. Thus, our design framework and analysis may be classified as a *thermal stress analysis* for real-time systems.

Categories and Subject Descriptors: C.2.2 [Real-Time Systems]: Control-Theoretic Thermal-Aware Systems Design

General Terms: Real-Time Systems, Thermal-Aware Systems, Multi-Mode Systems

Additional Key Words and Phrases: Control-theoretic systems, schedulability, EDF, Reactive systems, thermal resiliency, multi-mode system, thermal-aware system, thermal-aware periodic resource

ACM Reference Format:

Hettiarachchi, P., Fisher, N., Ahmed, M., Wang, L. Y., Wang, S., and Shi, W. 2013. A Design and Analysis Framework for Thermal-Resilient Hard-Real-Time Systems. *ACM Trans. Embedd. Comput. Syst.* 9, 4, Article 39 (March 2010), 25 pages.

DOI = 10.1145/0000000.0000000 <http://doi.acm.org/10.1145/0000000.0000000>

1. INTRODUCTION

Modern computer-controlled systems are often deployed in dynamic and unpredictable thermal operating environments. From the hardware-design perspective, material sci-

This research has been supported in part by the NSF (Grant Nos. CNS-0953585, CNS-1116787, CNS-1136007, and CNS-1205338), the Air Force Office of Scientific Research (Grant No. FA9550-10-1-0210), and two grants from Wayne State University's Office of Vice President of Research.

Author's addresses: P. Hettiarachchi and N. Fisher, M. Ahmed, S. Wang, and W. Shi, Department of Computer Science, Wayne State University; L. Wang, Department of Electrical and Computer Engineering, Wayne State University.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or permissions@acm.org.

© 2010 ACM 1539-9087/2010/03-ART39 \$15.00

DOI 10.1145/0000000.0000000 <http://doi.acm.org/10.1145/0000000.0000000>

entists and computer engineers use rigorous *thermal-stress analysis* techniques (e.g., see [Sergent and Krum 1998]) to determine how the underlying physical hardware will withstand applied internal and external thermodynamic forces. Unfortunately, equivalent analysis does not exist for determining the effects of (unpredictable) thermal stress on the performance of the systems software. While hardware capabilities such as dynamic power management (DPM) permit a computing system to reduce its power dissipation at run-time, many embedded systems have real-time constraints which may be adversely affected by unexpected changes in processor speed.

As an example of an embedded system where thermal-stress analysis is essential, consider microprocessors found in implantable medical devices (IMDs). IMDs are increasingly being used to treat various diseases and medical conditions (e.g., pacemakers for heart disease or neural implants to restore hearing/vision). However, recent studies [Kim et al. 2007; Lazzi 2005] have shown that the heat dissipated from IMDs due to the microprocessor activity is non-negligible. Thus, designing IMDs with minimum thermal dissipation is critical as medical research has shown that a temperature increase of even 1°C can have long-term effect on tissue [LaManna et al. 1989] and, in the extreme, death may even result from excessive tissue heating [Ruggera et al. 2003]. Complicating the safe thermal design of IMDs, body temperature naturally fluctuates over time and varies depending on location [Kelly 2006]. An IMD designer must balance (under temperature fluctuations) the real-time computational requirements of the device with the non-harmful thermal operating limits. In the presence of an increased surrounding temperature, an IMD will have to reduce its computational load to prevent tissue damage due to heat¹. However, as the correct and safe functioning of the IMD is an absolute requirement, the system designer requires techniques to formally verify the effect of different body temperatures on the correct operation of the IMD. Similarly, as a less safety-critical example, consider how the quality of audio/video decoding may degrade in a hand-held device as the system reacts to increases in temperature by reducing computational processing (e.g., via instruction fetch toggling). Ideally, a system designer would like to determine how much the performance will degrade under different thermal operating conditions.

Unfortunately, no current formal real-time design and analysis framework fully addresses the above setting. Recently-proposed control-theoretic frameworks exist for regulating processor temperature for *soft*-real-time systems (i.e., systems where jobs are permitted to “occasionally” miss computational deadlines) in an unpredictable thermal environment [Fu et al. 2010b,a]. While their results successfully show that it is possible to obtain stable and responsive thermal behavior and system utilization control, a system designer cannot use their approaches to *a priori* determine the amount of system-performance degradation due to changes in the thermal environment. Instead, the level of degradation can only be indirectly inferred via simulations of the system for different operating conditions. Furthermore, hard timing guarantees cannot be made in these frameworks. Techniques also already exist for permitting a trade-off between real-time QoS and processing resources (e.g., the *QoS-based resource allocation model* (GRAM) [Rajkumar et al. 1997]); however, while such techniques may guarantee real-time deadlines under a fixed level of resources, they cannot guarantee deadlines when a system must dynamically switch between real-time modes (due to the uncompleted execution remaining at mode transitions). Furthermore, none of these previously-proposed techniques can be used to obtain a precise, formal quantification of the thermal stress that the system can withstand.

¹As IMD microprocessors typically do not have DVS capabilities, an IMD may have to reduce non-essential tasks such as communication with other nodes in a body-area network [Timmons and Scanlon 2009].

In this paper, we address the challenge of determining the real-time guarantees in the presence of unpredictable dynamic environmental conditions. Towards this goal, we propose a framework and mechanisms for thermal-stress analysis in real-time systems. Our objective is *to develop techniques that permit a system designer to specify, a priori, a precise quantification of the hard-real-time performance degradation due to external thermal events, via a new system design metric called **real-time thermal resiliency***. Informally, real-time thermal resiliency is a prediction of the maximum external operating temperature at which a specified *real-time performance mode* (e.g., quality-of-service) may be guaranteed in the system steady-state (i.e., a time at which system properties have converged and do not change). To illustrate, consider a system with q different (system designer-defined) hard-real-time performance modes M_0, M_1, \dots, M_q where modes are ordered in increasing levels of real-time performance with M_q guaranteeing the highest level and M_0 the lowest. The real-time thermal resiliency of any mode M_i , denoted as $\Lambda(M_i, \mathcal{T}_{\text{ref}})$, is the predicted maximum external operating temperature for which the system will continue to operate (in the steady state) at performance mode M_i or higher and maintain a CPU reference temperature of \mathcal{T}_{ref} . Furthermore, if the external temperature exceeds $\Lambda(M_i, \mathcal{T}_{\text{ref}})$, then the system should automatically degrade to the next lowest performance mode M_{i-1} . The capability to define (at system-design time) thermal-resilient, real-time performance modes allows the system designer to specify how a system will gracefully and predictably degrade under external thermal stress; furthermore, the ability to accurately determine the real-time thermal resiliency of a performance mode provides a real-time system designer with a thermal-stress analysis framework analogous to stress analysis techniques in physical sciences and engineering. In the IMD example above, the thermal-resiliency function Λ may be used to determine (at design time) the body-temperature that a given set of tasks may safely operate at without doing damage to surrounding tissue.

§**Organization.** This paper presents a methodology for designing and analyzing thermal-resilient hard-real-time systems. Section 2 presents a high-level overview of our methodology and gives more detail on the contributions of this paper. Section 3 presents a brief review of previous work on thermal-aware (real-time and non-real-time) computer systems. Section 4 presents the hardware, real-time, and thermal models used throughout the paper. Section 5 details the design of our thermal-resilient controller. Section 6 derives thermal-resiliency function Λ for control system. Section 7 describes the results of our comparison with previous control systems via simulation and implementation upon testbed hardware. Our methodology provides formal system guarantees which require formal derivations and proofs. In the interest of space, we have deferred some of the formal proofs and derivations to a technical report available online [Hettiarachchi et al. 2011].

2. METHODOLOGY OVERVIEW

We now describe at a high level the major steps of our thermal-resilient design and analysis methodology.

- (1) **System Hardware Specification:** In the first step, the system designer must specify the processing and DPM capabilities of the system. Throughout this paper, we will be illustrating and validating our methodology upon an Intel Pentium IV 3.0 GHz single-core processor testbed. To match the rudimentary DPM capabilities often present in embedded processors, our testbed possesses the ability to only modulate the power modes of the system between active and inactive states. Section 4.1 gives more detail on the hardware model and our testbed implementation details.

- (2) **System Software Specification:** The system designer must specify the set of valid software modes M_0, M_1, \dots, M_q for the system. In Section 4.2, we discuss using the sporadic task model [Mok 1983] as a model for real-time workload of each software mode.
- (3) **Real-Time Mode Resource Allocation:** After the HW/SW specification steps, the designer must determine the minimum resource allocation under which the multi-mode system is schedulable. We discuss in Section 4.2 how recent techniques for schedulability analysis of hard-real-time systems where both the hardware and software change modes may be used in allocating sufficient processing time to each mode.
- (4) **Power/Thermal Model Evaluation:** Given the processing platform, we need an accurate power model in order to derive formal guarantees on the thermal resiliency of the system. Due to the duality between electrical and thermal circuits, we model the thermodynamics of our processing system using the resistance/capacitance (RC) circuits. We use *system identification* (SI) to identify the system parameters and evaluate the efficacy of our power-model choice. The details on the derived parameters for our hardware testbed are explained in the Section A.2 of the appendix.
- (5) **Control System Design:** We design a control structure based on optimal control theory. In this process, we use the SI parameters (determined in the previous step) to design the feedback gain parameters. We present details on our controller design in Section 5.
- (6) **System Simulation:** We build a system simulator which implements the real-time scheduling algorithm and control algorithm and simulates the real-time and thermal behavior of the system based on the resource allocations and power model derived in Steps 3 and 4. The details of our simulator are provided in Section 7.
- (7) **Thermal-Resiliency Function Calculation:** Given the real-time mode resource allocation, power model, controller, and simulator observations obtained from Steps 3, 4, 5, and 6 we can obtain a quantification of the thermal-resiliency function Λ . We give details on the derivation of this function in Section 6.
- (8) **System Validation:** We finally validate our system simulator and thermal-resiliency calculations in Section 7 by comparing directly with observations from our hardware testbed. Our comparison shows that the system simulator closely models the actual testbed behavior. Furthermore, we validate that our predicted thermal-resiliency Λ function is accurate by observing that it closely tracks the actual hardware testbed behavior.

While most of the steps above are standard practice in control system design, we would like to emphasize that our ability to ensure the hard-real-time schedulability of each mode in Step 3 and obtain *a priori* guarantees on thermal resiliency in Step 7 distinguishes our approach from previous thermal control for real-time systems.

3. RELATED WORK

In this section, we give a brief, high-level overview of previous research in both general (non-real-time), thermal-aware system design and real-time-specific thermal-aware design. For non-real-time systems, Brooks and Martonosi [2001] investigated major components of any dynamic thermal management scheme and suggested policies and mechanisms for implementing dynamic thermal management for current and future high-end CPUs. They evaluated the benefits of using dynamic thermal management to reduce the cooling system costs of CPUs and developed an architectural-level power modeling tool called Wattch. For the micro-architecture level of thermal modeling,

Skadron et al. [2003] proposed a compact, dynamic, and portable thermal model and a tool called *HotSpot* for use at the architecture level for micro-architectures.

For real-time systems in the online setting, Bansal and Pruhs [2005] explored algorithms for minimizing both peak-temperature and energy efficiency for online jobs with deadline constraints. In the off-line setting, previous work on scheduling under thermal constraints has followed two main approaches: reactive and proactive schedulers. In a reactive scheduler, the processor speed is reduced in response to a thermal trigger. Wang and Bettati [2008] studied schedulability analysis under the reactive setting. In the proactive setting, the speed schedule for the processor is determined at design time. Chen et al. [2009] addressed proactive scheduling for the periodic task model. Quan and Zhang [2009] consider feasibility analysis of leakage-aware periodic tasks under temperature constraints. However, previous work on both settings assumed either simple task models or the existence of “ideal” processor speeds. Our proposed control framework may be considered a proactive scheduler; however, we attempt to remove some ideal assumptions by working with only two power modes and the more general sporadic task model. Also, we consider the ambient temperature changes and analyze the effects on the task system due to its variation. Recent dynamic temperature management strategies also exist for multiprocessor real-time systems [Chen et al. 2007; Chantem et al. 2008; Fisher et al. 2009]; however, most of these focus upon static speed-assignment approaches and not a proactive schedule. Thermal analysis has also been studied in the context of web servers [Ferreira et al. 2007], but hard deadlines are not guaranteed. As mentioned in the introduction, work by Fu et al. [2010b] and Fu et al. [2010a] address handling unpredictable thermal events; however, the results do not provide any *a priori* guarantees that may be used to equate real-time performance and thermal resiliency.

4. MODELS

4.1. System Hardware Model and Testbed

For this paper, we consider a single processor system with rudimentary DPM capabilities of only *active* and *inactive* power modes. At any time $t > 0$, we denote the instantaneous CPU power as $\mathcal{P}_{\text{cpu}}(t)$. The processor dissipates thermal power at a constant rate $\mathcal{P}_{\text{cpu}}(t) = \mathcal{P}_{\text{act}}$ in the active mode and $\mathcal{P}_{\text{cpu}}(t) = \mathcal{P}_{\text{inc}}$ in the inactive mode. Also, we assume that processor consumes e_{act} amount of energy to activate from inactive mode and e_{inc} amount of energy to deactivate from the active mode. Even though the processor may be minimally active while in the low-power state, we will assume (as a pessimistic assumption for the purpose of schedulability analysis) that the processor is unavailable for task execution during this interval. If the aforementioned assumption does not hold, the system will behave “better” than the analysis and our results will continue to be valid. We believe this model of active/inactive modes is a very general model, applicable to a large number of available embedded processors with rudimentary DPM capabilities. For ideal processors with continuous power modes, $\mathcal{P}_{\text{cpu}}(t)$ may be selected from the range $[0, \mathcal{P}_{\text{act}}]$.

Our control system for the active/inactive processor will enforce strict periodic mode changes. For this purpose, we employ a recently proposed *thermal-aware periodic resource* [Ahmed et al. 2011] model, which is an extension of the well-known periodic resource model proposed by Shin and Lee [2008] for compositional real-time systems. In the thermal-aware periodic resource model, the processing resource is characterized with a two-tuple (Π, Θ) . The parameter Π is called the *resource period* and Θ is called the *resource capacity*. We will assume that Π is a non-negative integer (likely subject to the system tick granularity). The interpretation is that processor will be active for Θ amount of time at the beginning of each successive Π -length intervals. The ratio

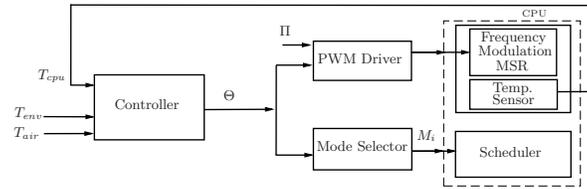


Fig. 1. The implementation details of the testbed. Note that the scheduler is responsible for EDF selection of jobs, activation of task threads to fill Θ and activation of idle thread during $\Pi - \Theta$, and thereby emulating the PWM cycle.

Θ/Π is called the *resource bandwidth*. Within each processor allocation, an arbitrary uniprocessor scheduling algorithm (e.g., EDF or RM) may be employed to schedule the underlying task system (see next subsection). See Figure 2 for an illustration of the thermal-aware periodic resource.

As a case study of our methodology, we have built a hardware testbed using an Intel Pentium IV 3.0 GHz single core processor running a modified Linux kernel (2.6.33.7.2-rt30 PREEMPT_RT). The low power CPU on our testbench does not have a System Developer Interface to measure the on-die temperature directly. We follow the procedure given in the Intel Documentation [Int 2000] and install a T-type thermocouple on the CPU die. We use Phidgets 4-port temperature sensor board to measure the environment, air, and the on-die temperature through the USB driver and allows us to directly interface the sensors with the testbed software.

We develop a loadable kernel module to activate and vary the frequency modulation level at run-time. We use Model Specific Registers (MSR) to control the frequency modulation ratio in the clock and select the higher and the lowest frequency modulation indices to emulate the low and the higher power levels. We use 12.5% and 87.5% modulation ratios in the `IA32_CLOCK_MODULATION` MSR for active and inactive power mode emulation.

We develop a multi-threaded application using Linux native posix thread libraries (NTPL). Our application consists of a scheduler simulator and a thread activator where the schedule simulator selects the EDF based jobs from the local ready-queue and dispatches them into a thread activator. The thread activator consists of a very high priority thread (priority is set to higher than the threaded IRQ handlers), emulates the schedule tick in the Linux kernel in higher level abstraction. Similar to the Linux kernel scheduler tick, the thread activator sleeps until it wakes up accurately in the scheduling boundaries. Our thread activator wakes up in unequal tick intervals to schedule jobs, raises the appropriate thread which should have the priority, and goes back to the sleeps. The jobs are selected by the schedule simulator according to EDF. This process repeats and the amount of time allocates to each job depends on EDF and the total time depends on the Θ given by the optimal controller.

Figure 1 provides a high-level overview of the workflow for the different components of our framework. The controller after sampling the temperatures determines the capacity. The capacity is given to the PWM controller and the real-time performance mode selector. The PWM modulates the frequency of the CPU via the MSR and an OS Scheduler (EDF) determines how to schedule the selected performance mode within the PWM duty cycle. Our temperature sensors sample the temperatures and the process iterates ad infinitum.

4.2. System Software Model

In the introduction, we proposed a system model of *real-time performance modes* M_1, \dots, M_q . For the purpose of this paper, we will assume each performance mode

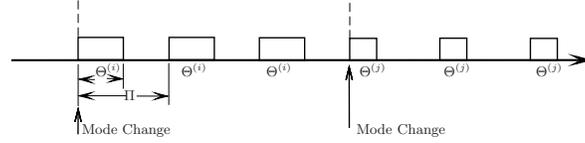


Fig. 2. The sampling and mode change in our thermal control system. The blocks indicate time periods during which the processor is active under the thermal-aware periodic resource model. Sporadic tasks are scheduled within the activation blocks.

M_i is characterized by a *sporadic task system*² [Mok 1983] with n_i tasks and the resource capacity $\Theta^{(i)}$. That is, $M_i = \left(\{ \tau_1^{(i)}, \tau_2^{(i)}, \dots, \tau_{n_i}^{(i)} \}, \Theta^{(i)} \right)$ where each $\tau_j^{(i)} \in M_i$ is a sporadic task characterized by a three-tuple $(e_j^{(i)}, d_j^{(i)}, p_j^{(i)})$ and $\Theta^{(i)}$ is the minimum capacity required to meet the deadlines of the tasks of M_i . (Note that we are abusing notation by allowing M_i to represent the set of tasks and the two-tuple of the mode's task system and required resource capacity.) In this three-tuple representation for a task, $e_j^{(i)}$ is the *worst-case execution requirement*, $d_j^{(i)}$ is the *relative deadline*, and $p_j^{(i)}$ is the *minimum inter-arrival separation parameter* (historically called the “period”). A sporadic task $\tau_j^{(i)}$ may produce a (potentially infinite) sequence of jobs, where each job has an execution requirement of $e_j^{(i)}$ time units and must complete $d_j^{(i)}$ time units after its arrival. The first job of $\tau_j^{(i)}$ may arrive at any time after system-start time; however, successive jobs of $\tau_j^{(i)}$ must arrive at least $p_j^{(i)}$ time units apart. For this paper, we assume that the resource period Π is identical in all modes. For mode M_i , a resource capacity of $\Theta^{(i)}$ is provided every resource period. Figure 2 illustrates the processing-time allocation in two different modes.

We will assume that there is an ordering of real-time performance modes based on their “computational requirements” to meet all of a mode’s deadlines. The relation $M_i \succeq M_j$ indicates that M_i is more computationally intensive than M_j . For notational convenience, we will assume that mode M_0 represents the mode with no tasks and $\Theta^{(0)}$ equal to zero. Furthermore, for this paper, we assume that the modes are well-ordered and have been indexed in increasing order of computational requirements; i.e., $M_0 \preceq M_1 \preceq M_2 \preceq \dots \preceq M_q$. While there are many possible ways to define the \preceq relation, the only ordering required from the perspective of our thermal control is that $M_i \preceq M_j$, if and only if, $\Theta^{(i)} \leq \Theta^{(j)}$; i.e., to reduce the temperature of the system, we need to decrease the processing-time allocation.

Our model does not require any particular mode-change semantics to be adopted. Some potential options for dealing with incompletely-executed jobs upon a mode change are: (i) aborting any incomplete jobs; (ii) delaying the release of jobs in the new mode until all jobs of the old mode have completed; and (iii) allowing jobs of the new mode to be released, as soon as legally allowable, while jobs of the old mode are still active. For the purposes of our hardware testbed and simulations (Section 7), we assume option (iii).

The scheduling of real-time performance mode M_i upon the thermal-aware periodic resource may be done by any uniprocessor real-time scheduling algorithm (e.g., earliest-deadline-first or rate-monotonic [Liu and Layland 1973]). However, $\Theta^{(i)}$ must be sufficiently large for the scheduling algorithm to correctly schedule all jobs of the task set of M_i (i.e., $\{ \tau_1^{(i)}, \tau_2^{(i)}, \dots, \tau_{n_i}^{(i)} \}$) and (potentially) any jobs from the previous mode that have not completed by the mode change. To obtain a proper resource allocation

²Note, we will be assuming the sporadic task model throughout our objectives, but the results could be extended to other task models without much change.

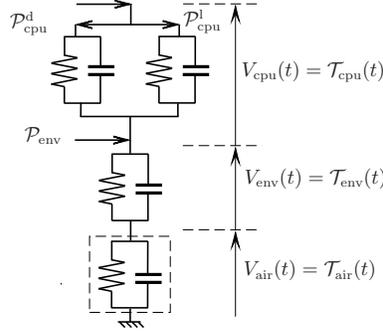


Fig. 3. The basic equivalent circuit for a working CPU and its working environment

tion, $\Theta^{(i)}$, for each mode, we use our recently-developed hard-real-time schedulability test (for EDF scheduling under hardware/software mode changes in the periodic resource model) to search for a safe value of $\Theta^{(i)}$ for each mode [Fisher and Ahmed 2011] to ensure that deadlines are always met. The multi-modal schedulability analysis ensures that for any valid sequence of mode changes and valid set of job arrivals under the sporadic task model that the EDF scheduler will always meet all deadlines. The analysis works by determining the maximum workload carried from one mode to another and testing whether this “carry-in” will cause a deadline miss.

4.3. Power/Thermal Model

We use the duality principle in electrical and thermal circuits to describe the dynamics of the power dissipating source using electrical resistance/capacitance (RC) circuits. Figure 3 shows the basic equivalent circuit for the CPU and its surrounding environment. We assume that total dissipated power of the CPU \mathcal{P}_{cpu} is equal to the sum of the power due to dynamic current $\mathcal{P}_{\text{cpu}}^d$ and power due to leakage current $\mathcal{P}_{\text{cpu}}^l$. Furthermore, we assume that the temperature-dependant leakage power may be closely approximated by a linear function of CPU temperature [Liu et al. 2007].

Let $V_{\text{cpu}}(t)$, $V_{\text{env}}(t)$, and $V_{\text{air}}(t)$ represent the equivalent voltages for temperatures of the CPU, environment, and air (room) respectively. Let \mathcal{T}_{cpu} be the instantaneous relative temperature of the CPU with respect to the immediate environment (e.g., CPU casing), \mathcal{T}_{env} be the relative temperature of the immediate environment with respect to the room air temperature, and \mathcal{T}_{air} be the (absolute) room air temperature. For example, if \mathcal{T}_{air} is 20°C , \mathcal{T}_{env} is 10°C , and \mathcal{T}_{cpu} is 15°C , then the absolute temperature of the CPU is 45°C .

Let $\mathcal{P}_{\text{cpu}}^d(t)$, $\mathcal{P}_{\text{cpu}}^l(t)$, and $\mathcal{P}_{\text{env}}(t)$ represent, respectively, the dynamic CPU, leakage CPU, and environment power dissipation. Let R_{cpu}^d , R_{cpu}^l , R_{env} , C_{cpu}^d , C_{cpu}^l , and C_{env} represent the dynamic and leakage thermal resistance, environment resistance, CPU dynamic and leakage capacitance, and environment capacitance. Finally, let $\sigma_1 \stackrel{\text{def}}{=} \frac{1}{C_{\text{cpu}}^d + C_{\text{cpu}}^l}$ and k_T and k_C represent processor-dependent constants used in approximating the temperature-dependant leakage current. Applying Kirchhoff’s circuit

5.1. State-Space Model Basics

We use the standard state-space model to represent continuous-time (ideal) system

$$\begin{aligned}\dot{x}(t) &= Ax(t) + Bu(t) + f, \\ y(t) &= Cx(t),\end{aligned}\tag{8}$$

where $x(t)$, $u(t)$, and $y(t)$ represent the state vector, the input vector, and the output vector, respectively. A , B , and C represent the system matrices and f represents a constant vector. Both the state matrices and constant vector are time-invariant quantities.

Since we have a computer-controlled discrete-time system, we will use following state-space mode for the discrete-time controller for active/inactive modes. For a sampling interval T_s , $u(t)$ is a constant and the sampled system of Equation (8) is

$$\begin{aligned}x((k+1)T_s) &= Gx(kT_s) + Hu(kT_s) + \tilde{f}, \\ y(kT_s) &= \tilde{C}x(kT_s),\end{aligned}\tag{9}$$

where $G = e^{AT_s}$, $H = \int_0^{T_s} e^{At} B dt$, $\tilde{C} = C$, and $\tilde{f} = \int_0^{T_s} e^{At} f dt$. The term e^{At} can be computed by $\mathcal{L}^{-1}\{(sI - A)^{-1}\}$, where \mathcal{L}^{-1} is the inverse Laplace transform. In the remainder of the document, we abuse the notation by representing $x(kT_s)$ as $x(k)$, $x((k+1)T_s)$ as $x(k+1)$, $u(kT_s)$ as $u(k)$, and $y(kT_s)$ as $y(k)$. The above definitions may be found in any textbook on discrete-time control theory [Ogata 1995].

5.2. Continuous Power Modes

As a first step towards our goal of designing a control-theoretic framework for thermal stress analysis, we employ *linear quadratic (LQ) optimal control* for real-time thermal management. Our design consists of an optimal state feedback and an integrator that regulates the dynamics of the system. An LQ controller enables us to design an efficient and low-overhead controller, derive the feedback parameters before runtime (used in thermal-resiliency analysis), and smoothly track our reference input. In the future, we plan on applying more complex and robust controllers (e.g., \mathcal{H}_∞ controllers) to decrease the controller's sensitivity to modeling inaccuracy and noise. However, as observed in the simulations and experiments of Section 7, our current LQ design is appropriately responsive to changes in environmental temperature.

In our system model, we specify the thermal power of the CPU as the control to the system. The controller is designed to follow the temperature reference, T_{ref} . In our design, we consider $\mathcal{T}_{\text{cpu}}(t)$ as one of the variable to be controlled and $\mathcal{P}_{\text{cpu}}^{\text{d}}(t)$ as a manipulated variable (equivalent to $y(t)$ and $u(t)$, respectively, in continuous state-space model). The basic control structure is given in Figure 4.

From Equations (3) and (5), the continuous-time state space model can be written as

$$\begin{bmatrix} \dot{\mathcal{T}}_{\text{cpu}}(t) \\ \dot{\mathcal{T}}_{\text{env}}(t) \end{bmatrix} = \begin{bmatrix} -\beta_1 & k_T \sigma_1 \\ k_T \sigma_2 & -\beta_2 \end{bmatrix} \begin{bmatrix} \mathcal{T}_{\text{cpu}}(t) \\ \mathcal{T}_{\text{env}}(t) \end{bmatrix} + \begin{bmatrix} \sigma_1 \\ \sigma_2 \end{bmatrix} \mathcal{P}_{\text{cpu}}^{\text{d}}(t) + \begin{bmatrix} 0 \\ \sigma_2 \end{bmatrix} \mathcal{P}_{\text{env}}(t).\tag{10}$$

While our analysis below is in the continuous-time domain, a discrete-time control system approach would be applied in an actual computer implementation. Therefore, we now note that we may easily convert the continuous-state space model to the discrete-time sampled system, $x(k+1) = Gx(k) + Hu(k) + f$ from the continuous-time state matrices $A = \begin{bmatrix} -\beta_1 & k_T \sigma_1 \\ k_T \sigma_2 & -\beta_2 \end{bmatrix}$ and $B = \begin{bmatrix} \sigma_1 \\ \sigma_2 \end{bmatrix}$ where k is the sampling index, T_s is sampling interval, and G and H can be calculated as described in Section 5.1. For

our given system, $x(k) \equiv \begin{bmatrix} \mathcal{T}_{\text{cpu}}(k) \\ \mathcal{T}_{\text{env}}(k) \end{bmatrix}$ and $u(k) \equiv [\mathcal{P}_{\text{cpu}}(k)]$ where we are again abusing notation for the \mathcal{T} and \mathcal{P} functions.

To eliminate steady state tracking error, we design our control system with an integrator. Define an additional error vector $v_e(t)$ in continuous time as,

$$\begin{aligned} v_e(t) &\stackrel{\text{def}}{=} \int_0^t (\mathcal{T}_{\text{ref}} - \mathcal{T}(t) - \mathcal{T}_{\text{air}}(t)) dt \\ \dot{v}_e(t) &\stackrel{\text{def}}{=} \mathcal{T}_{\text{ref}} - \mathcal{T}_{\text{air}}(t) - \mathcal{T}(t) \\ &= -C \begin{bmatrix} \mathcal{T}_{\text{cpu}}(t) \\ \mathcal{T}_{\text{env}}(t) \end{bmatrix} + \mathcal{T}_{\text{ref}} - \mathcal{T}_{\text{air}}(t) \end{aligned} \quad (11)$$

where $C = [1, 1]$. Then, the system input is calculated with a gain $K_o = [\gamma_1, \gamma_2]$ and integral constant γ_I in the following equation.

$$\begin{aligned} \mathcal{P}_{\text{cpu}}^{\text{d}}(t) &= -K_o \begin{bmatrix} \mathcal{T}_{\text{cpu}}(t) \\ \mathcal{T}_{\text{env}}(t) \end{bmatrix} + \gamma_I v_e(t) \\ &= -\left((\gamma_1) \mathcal{T}_{\text{cpu}}(t) + (\gamma_2) \mathcal{T}_{\text{env}}(t) \right) + \gamma_I \int_0^t (\mathcal{T}_{\text{ref}} - \mathcal{T}_{\text{air}}(t) - \mathcal{T}_{\text{cpu}}(t) - \mathcal{T}_{\text{env}}(t)) dt. \end{aligned} \quad (12)$$

We employ standard techniques from optimal control theory to derive K_o and γ_I and prove stability. In our derivation of system stability, we use the following two results which can be found in any standard text on control theory [Dorf and Bishop 2000; Nise 2000; Ogata 1995].

LEMMA 1 (FROM [DORF AND BISHOP 2000]). *The system of Equation (8) is completely controllable if there exists an unconstrained $u(t)$ such that it can control any initial state $x(t_0)$ to any desired final state x_f in a finite time, $t_0 \leq t \leq T$. The property of completely controllable can be determined by examining the algebraic condition*

$$\text{rank}[B \quad AB \quad A^2B \quad \dots \quad A^{m-1}B] = m, \quad (13)$$

where, A is $m \times m$ and B is $m \times r$ matrix.

LEMMA 2 (FROM [OGATA 1995]). *A discrete-time linear time invariant (LTI) system is asymptotically stable if and only if its all eigenvalues of G lie inside the unit circle.*

We derive the augmented model that is used to obtain the optimality of the system. Consider an instance where system is completely stable and has reached steady state. We denote the input, states, and the integrator error (described in Equation (11)) of this special instance of the system by $\mathcal{P}_{\text{cpu}}(t_\infty)$, $\mathcal{T}_{\text{cpu}}(t_\infty)$, $\mathcal{T}_{\text{env}}(t_\infty)$ and $v_e(t)$ respectively. Therefore, from the Equation (10) we get,

$$\begin{bmatrix} \dot{\mathcal{T}}_{\text{cpu}}(t) - \dot{\mathcal{T}}_{\text{cpu}}(t_\infty) \\ \dot{\mathcal{T}}_{\text{env}}(t) - \dot{\mathcal{T}}_{\text{env}}(t_\infty) \end{bmatrix} = \begin{bmatrix} -\beta_1 & k_T \sigma_1 \\ k_T \sigma_2 & -\beta_2 \end{bmatrix} \begin{bmatrix} \mathcal{T}_{\text{cpu}}(t) - \mathcal{T}_{\text{cpu}}(t_\infty) \\ \mathcal{T}_{\text{env}}(t) - \mathcal{T}_{\text{env}}(t_\infty) \end{bmatrix} + \begin{bmatrix} \sigma_1 \\ \sigma_2 \end{bmatrix} (\mathcal{P}_{\text{cpu}}^{\text{d}}(t) - \mathcal{P}_{\text{cpu}}^{\text{d}}(t_\infty)). \quad (14)$$

Also, from the Equation (11), we get,

$$\dot{v}_e(t) - \dot{v}_e(t_\infty) = -C \begin{bmatrix} \mathcal{T}_{\text{cpu}}(t) - \mathcal{T}_{\text{cpu}}(t_\infty) \\ \mathcal{T}_{\text{env}}(t) - \mathcal{T}_{\text{env}}(t_\infty) \end{bmatrix}. \quad (15)$$

Now, combining the Equation (14) and (15), we define our higher order system as,

$$\dot{e}(t) = \hat{A}e(t) + \hat{B}u_e(t), \quad (16)$$

where, $e(t) = \begin{bmatrix} \mathcal{T}_{\text{cpu}}(t) - \mathcal{T}_{\text{cpu}}(t_\infty) \\ \mathcal{T}_{\text{env}}(t) - \mathcal{T}_{\text{env}}(t_\infty) \\ v_e(t) - v_e(t_\infty) \end{bmatrix}$, $u_e(t) = \mathcal{P}_{\text{cpu}}^d(t) - \mathcal{P}_{\text{cpu}}^d(t_\infty)$, $\hat{A} = \begin{bmatrix} A & 0 \\ -C & 0 \end{bmatrix}$, and $\hat{B} = [B \ 0]^T$.

We select the feedback gain $\hat{\gamma}$ such that,

$$u_e(t) = -\hat{K}e(t), \quad (17)$$

where, $\hat{K} = [K_o \ -\gamma_I]$. The above state-space and the control gain parameters are valid for a continuous-time controller. So, we may obtain the discrete-time state-space matrices for the augmented model (i.e, G and H) from \hat{A} and \hat{B} via the transformation described after Equation (9). In LQ optimal control, the objective is to design the controller to minimize some performance index. A standard LQ performance index is given by

$$J \stackrel{\text{def}}{=} \frac{1}{2} \sum_{k=0}^{\infty} (e(k)^T Q e(k) + u_e^T(k) R u_e(k)), \quad (18)$$

where Q and R are arbitrary symmetric matrices of size $m \times m$ and $r \times r$ such that $Q \geq 0$ (positive semi definite), $R > 0$ (positive definite). (In our system given in Equation (10), m is two and r is one). It is easy to show that for a Linear Time Invariant (LTI) system, (Refer to Ogata [1995]), the optimal state feedback can be obtained as,

$$u_e(k) = -\hat{K}e(k), \quad (19)$$

where \hat{K} is the feedback gain defined as

$$\hat{K} = (R + H^T P H)^{-1} H^T P G, \quad (20)$$

and where P is the positive definite solution of the algebraic Riccati equation below,

$$P = Q + G^T P G - G^T P H (R + H^T P H)^{-1} H^T P G.$$

From the above, it may be shown [Ogata 1995] that the optimal performance index can be calculated as $J_{\min} = \frac{1}{2} e^T(0) P e(0)$.

It is well known [Ogata 1995] that the feedback control (i.e., \hat{K}) results in an asymptotically stable closed-loop system according to Lemma 2. Obviously, stable choices of K_o and γ_I for the original (non-augmented) system can be immediately obtained from the derived \hat{K} .

5.3. Active/Inactive Power Modes

Since the CPU power cannot be varied continuously, the controller designed in the previous section cannot be directly applied to the setting of discrete active/inactive power modes. In this section, we extend the design of the continuous power modes controller described in the previous section to the active/inactive power mode setting by applying pulse-width modulation (PWM) techniques. Recall in Section 4 that we stated the active/inactive power modes will be modeled via the thermal-aware periodic resource model with parameters Π and Θ . Thus, to control the system via this model, we must choose the appropriate values of Π and Θ . The Π value is a design parameter which may be chosen at controller design-time and will be assumed fixed throughout controller execution. Typically, a smaller value of Π will increase the system schedulability; however, a larger value of Π will decrease the overhead potentially incurred by switching between the active and inactive power modes. (See Ahmed et al. [2011] for algorithms for determining Π in the thermal setting). The only constraint that our framework places on the chosen value of Π is that it must evenly divide the sampling interval length T_s (i.e., $T_s = \kappa \Pi$ for some $\kappa \in \mathbb{N}^+$).

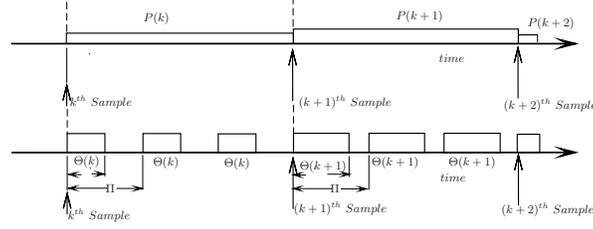


Fig. 5. The simplified power and modulation relationship

Since we have only two power modes, we cannot arbitrarily set the power level. However, we may change the assigned resource capacity between sampling periods to approximate arbitrary power levels. Therefore, the assigned resource capacity will be the manipulated variable in our PWM system. The periodic resource capacity ($\Theta(k)$) and resource period (Π) can be respectively viewed as the pulse duration and duty cycle of the PWM. Let $\Theta(k)$ denote the value of the resource capacity over the k 'th sampling period. For determining the $\Theta(k)$ value, we use a method based on the *principle of equivalent areas* (PEA) for converting any arbitrary input signal into an equivalent PWM signal [Gel'ig and Churilov 1998]. First, note that in a discrete-time system using *zero-order hold* (ZOH), the input signal is held constant over the sampling period. Specifically, for the k 'th sampling interval, the input $\mathcal{P}_{\text{cpu}}^{\text{d}}(k)$ is held over the T_s -length interval, resulting in a total energy dissipation of $T_s \cdot \mathcal{P}_{\text{cpu}}^{\text{d}}(k)$ over the interval. To get the equivalent area (i.e., energy) as the (ideal) system with continuous power modes, we must set $\Theta(k)$ such that the periodic modulations between the power modes of \mathcal{P}_{act} and \mathcal{P}_{inc} dissipate the equivalent amount of energy over the T_s -length interval. Figure 5 illustrates the area equivalence between the continuous and PWM controllers. It is easy to see that a smaller T_s gives a better PWM approximation. However, our controller needs to follow a system with relatively slower (thermal) dynamics. Thus, for efficiency, we select relatively larger T_s and higher κ value. Also, even under varying air and environmental conditions, the resource capacity, Θ does not change rapidly due to slower system dynamic. Therefore, the same mode will continue to hold over several sampling periods before change occurs. Furthermore, in the steady state (when the environment or air temperature does not change much), the system will change modes very infrequently. The technical report version of this paper [Hettiarachchi et al. 2011] describes how to choose T_s to minimize error due to the PWM approximation.

More formally, we may derive the following relationship between $\mathcal{P}_{\text{cpu}}^{\text{d}}(k)$ and $\Theta(k)$,

$$\begin{aligned} \kappa \Pi \mathcal{P}_{\text{cpu}}^{\text{d}}(k) &= \kappa \left(e_{\text{act}} + \int_0^{\Theta(k)} \mathcal{P}_{\text{act}} dt + e_{\text{inc}} + \int_{\Theta(k)}^{\Pi} \mathcal{P}_{\text{inc}} dt \right) \\ \Rightarrow \mathcal{P}_{\text{cpu}}^{\text{d}}(k) &= \left(\frac{\mathcal{P}_{\text{act}} - \mathcal{P}_{\text{inc}}}{\Pi} \right) \Theta(k) + \mathcal{P}_{\text{inc}} + \frac{1}{\Pi} (e_{\text{act}} + e_{\text{inc}}). \end{aligned} \quad (21)$$

The PWM controller pseudocode is presented in Algorithm 1. The controller proposed here consists of two integrated operations: the thermal controller and the PWM modulator. The first step is to obtain the CPU temperature at t_ℓ (Line 2 of Algorithm 1). The error is then calculated by taking the difference between the reference temperature and the CPU temperature (Line 3). The error is integrated into the error vector and added to vector sum of the integrated error in the next line (Line 4). After which, the power input is calculated (Line 5) and the equivalent Θ is calculated from the property of Equation (21) (Line 6). Finally, the appropriate mode is selected (Line 7), the

Algorithm 1 Control Algorithm

Require: Reference Temperature \mathcal{T}_{ref} ; Feedback Gain $K \equiv [\gamma_1, \gamma_2]$; Integral Constant γ_I ; PWM Period Π ; Number of PWM periods in a sampling period κ .

- 1: **while** At beginning of sampling period $[t_\ell, t_{\ell+1}) : t_\ell \equiv \kappa\ell\Pi$ **do**
- 2: **Sample** $\mathcal{T}_{\text{cpu}}(t_\ell) + \mathcal{T}_{\text{env}}(t_\ell) + \mathcal{T}_{\text{air}}(t_\ell)$.
- 3: $\dot{v}_e(t_\ell) = \mathcal{T}_{\text{ref}} - (\mathcal{T}_{\text{cpu}}(t_\ell) + \mathcal{T}_{\text{env}}(t_\ell) + \mathcal{T}_{\text{air}}(t_\ell))$
- 4: $Tot_v_e(t_\ell) = Tot_v_e(t_{\ell-1}) + \gamma_I \kappa \Pi \frac{(\dot{v}_e(t_\ell) + \dot{v}_e(t_{\ell-1}))}{2}$
- 5: $\mathcal{P}_{\text{cpu}}(t_\ell) = \left(Tot_v_e(t_\ell) - (\gamma_1 \mathcal{T}_{\text{cpu}}(t_\ell) + \gamma_2 \mathcal{T}_{\text{env}}(t_\ell)) \right)$
- 6: $\Theta(t_\ell) = \min \left(\Pi \times \frac{(\mathcal{P}_{\text{cpu}}(t_\ell) - \mathcal{P}_{\text{inc}})}{\mathcal{P}_{\text{act}} - \mathcal{P}_{\text{inc}}}, \Pi \right)$
- 7: $i = \max\{j \in \mathbb{Z}_{q+1} \mid \Theta^{(j)} \leq \Theta(t_\ell)\}$
- 8: **Update** real-time performance mode to M_i .
- 9: **Set** PWM to operate at period of Π and width of $\Theta(t_\ell)$.
- 10: **end while**

mode change is performed (Line 8), and the pulse-width modulator is invoked for the next $\kappa \Pi$ -length intervals (Line 9). It is important to note that $\Theta(t_\ell)$ calculated in Line 6 does not have to be equal the $\Theta^{(j)}$ for the selected mode; we must only select the highest mode with $\Theta^{(j)} \leq \Theta(t_\ell)$. (If $\Theta(t_\ell)$ is larger, we are only giving the mode more processing than it requires.) It should also be observed that all operations, except for finding the appropriate mode, may be done in $O(1)$ time. Finding the highest real-time performance mode that may execute can be done in $O(\lg q)$ time (via binary search) where q is the number of real-time performance modes.

6. THERMAL-RESILIENCY CALCULATION

In this section, we explain how to derive the real-time thermal resiliency $\Lambda(M_i, \mathcal{T}_{\text{ref}})$ for a given real-time performance mode M_i and reference temperature \mathcal{T}_{ref} . Assuming a steady-state error of zero, we will now briefly outline how to obtain a solution for $\Lambda(M_i, \mathcal{T}_{\text{ref}})$.⁴ Assume that we have reached the steady-state by the $(k-1)$ 'th sampling period. Therefore, $\mathcal{T}_{\text{cpu}}(k) = \mathcal{T}_{\text{cpu}}(k-1)$, $\mathcal{T}_{\text{env}}(k) = \mathcal{T}_{\text{env}}(k-1)$, $\mathcal{T}_{\text{air}}(k) = \mathcal{T}_{\text{air}}(k-1)$, and $\Theta(k) = \Theta(k-1)$. Substituting the temperature equalities into Equations (6) and (7) allows us to solve for $\mathcal{T}_{\text{cpu}}(k)$ and $\mathcal{T}_{\text{env}}(k)$ to obtain a function of $\mathcal{T}_{\text{air}}(k)$, \mathcal{T}_{ref} , and $\Theta(k)$. Since we are interested in obtaining $\Lambda(M_i, \mathcal{T}_{\text{ref}})$, we may fix \mathcal{T}_{ref} and $\Theta(k) = \Theta^{(i)}$. Since the steady-state error is zero, we also have

$$\mathcal{T}_{\text{ref}} = \mathcal{T}_{\text{cpu}}(k) + \mathcal{T}_{\text{env}}(k) + \mathcal{T}_{\text{air}}(k). \quad (22)$$

Combining Equation 22 with the function of $\mathcal{T}_{\text{air}}(k)$ obtained from $\mathcal{T}_{\text{cpu}}(k)$ and $\mathcal{T}_{\text{env}}(k)$ allows us to solve for $\mathcal{T}_{\text{air}}(k)$. Thus, solving the entire system results in a value for $\mathcal{T}_{\text{env}}(k) + \mathcal{T}_{\text{air}}(k)$ (i.e., value of $\Lambda(M_i, \mathcal{T}_{\text{ref}})$). The resulting expression is quite complicated as it requires solutions to second-order inhomogeneous equations.

We first calculate the \mathcal{T}_{cpu} as follows,

$$\begin{aligned} \mathcal{T}_{\text{cpu}}((\zeta\kappa + \kappa)\Pi) &= \mathcal{T}_{\text{cpu}}(\zeta\kappa\Pi) + \sum_{i=0}^{\kappa-1} \sum_{j=1}^2 \left(\mathcal{C}_{(j)_{\text{inc}}}((\zeta\kappa + i)\Pi + \Theta)(e^{r(j)(\Pi - \Theta)} - 1) \right. \\ &\quad \left. + \mathcal{C}_{(j)_{\text{act}}}((\zeta\kappa + i)\Pi)(e^{r(j)\Theta} - 1) \right). \end{aligned} \quad (23)$$

⁴The approach may be generalized when there is bounded steady-state error. However, the approach will be similar, and we omit the details due to space.

At the stability, $\mathcal{T}_{\text{cpu}}(\zeta\kappa\Pi)$ stays at a steady value and therefore, $\mathcal{T}_{\text{cpu}}((\zeta\kappa + \kappa)\Pi)$ and $\mathcal{T}_{\text{cpu}}(\zeta\kappa\Pi)$ are the same. Further, if the CPU does not vary the temperature within a single sampling period, the CPU should maintain the same temperature in each resource period Π intervals (For same Θ , same \mathcal{T}_{cpu} and \mathcal{T}_{env} temperature at successive stages). Therefore, we consider the CPU temperature for two adjacent resource periods and conclude,

$$\sum_{j=1}^2 \left(\mathcal{C}_{(j)_{inc}}(\zeta\kappa\Pi + \Theta)(e^{r^{(j)}(\Pi-\Theta)} - 1) + \mathcal{C}_{(j)_{act}}(\zeta\kappa\Pi)(e^{r^{(j)}\Theta} - 1) \right) = 0, \quad (24)$$

because, $\mathcal{T}_{\text{cpu}}^{\text{inc}}((\zeta\kappa + 1)\Pi) = \mathcal{T}_{\text{cpu}}^{\text{act}}(\zeta\kappa\Pi)$ as per to the above argument. Then we further simplify the Equation (24) as follows (details in the tech report),

$$\begin{aligned} &\Rightarrow \mathcal{T}_{\text{env}}(\zeta\kappa\Pi + \Theta) \left(\mathcal{P}_4(\Theta) \right) + \mathcal{T}_{\text{cpu}}(\zeta\kappa\Pi + \Theta) \left(\mathcal{P}_3(\Theta) \right) + \mathcal{T}_{\text{cpu}}(\zeta\kappa\Pi) \left(\mathcal{P}_1(\Theta) \right) \\ &+ \mathcal{T}_{\text{env}}(\zeta\kappa\Pi) \left(\mathcal{P}_2(\Theta) \right) + \mathcal{P}_A(\Theta) = 0 \end{aligned} \quad (25)$$

where, $\mathcal{P}_4(\Theta) = (\mathcal{G}_4(e^{r_1(\Pi-\Theta)} - 1) + \mathcal{G}_8(e^{r_2(\Pi-\Theta)} - 1))$, $\mathcal{P}_3(\Theta) = \mathcal{T}_{\text{cpu}}(\zeta\kappa\Pi + \Theta) (-\mathcal{G}_3(e^{r_1(\Pi-\Theta)} - 1) - \mathcal{G}_7(e^{r_2(\Pi-\Theta)} - 1))$, $\mathcal{P}_1(\Theta) = \mathcal{T}_{\text{cpu}}(\zeta\kappa\Pi) (-\mathcal{G}_1(e^{r_1\Theta} - 1) - \mathcal{G}_5(e^{r_2\Theta} - 1))$, $\mathcal{P}_2(\Theta) = \mathcal{T}_{\text{env}}(\zeta\kappa\Pi) (\mathcal{G}_2(e^{r_1(\Pi-\Theta)} - 1) + \mathcal{G}_6(e^{r_2\Theta} - 1))$, and $\mathcal{P}_A(\Theta) = \mathcal{G}_A(e^{r_1\Theta} - 1) + \mathcal{G}_B(e^{r_1(\Pi-\Theta)} - 1) + \mathcal{G}_C(e^{r_2\Theta} - 1) + \mathcal{G}_D(e^{r_2(\Pi-\Theta)} - 1)$.

In Equation (25), we use the definitions of \mathcal{C} for $\zeta\kappa\Pi$ and $(\zeta\kappa\Pi + \Theta)$ time instances as shown below for $i = 1, 2$. Let \bar{i} equal two if i equals one and one if i equals two.

$$\begin{aligned} \mathcal{C}_{i_{act}}(\zeta\kappa\Pi) &= \frac{(-1)^{(i)}}{r_2 - r_1} \left(\begin{array}{l} r_{\bar{i}}\mathcal{C}_{3_{inc}}(\zeta\kappa\Pi) \\ +\sigma_1(\mathcal{P}_{act} + k_C + k_T\mathcal{T}_{\text{env}}(\zeta\kappa\Pi)) \\ -(\beta_1 + r_{\bar{i}})\mathcal{T}_{\text{cpu}}(\zeta\kappa\Pi) \end{array} \right) \\ &= \frac{(-1)^{(i)}}{r_2 - r_1} (\mathcal{G}_{A_i} + \mathcal{G}_B\mathcal{T}_{\text{env}}(\zeta\kappa\Pi) - \mathcal{G}_{C_i}\mathcal{T}_{\text{cpu}}(\zeta\kappa\Pi)), \\ \mathcal{C}_{i_{inc}}(\zeta\kappa\Pi + \Theta) &= \frac{(-1)^i}{r_2 - r_1} \left(\begin{array}{l} r_{\bar{i}}\mathcal{C}_{3_{inc}}(\zeta\kappa\Pi + \Theta) \\ +\sigma_1(\mathcal{P}_{inc} + k_C + k_T\mathcal{T}_{\text{env}}(\zeta\kappa\Pi + \Theta)) \\ -(\beta_1 + r_{\bar{i}})\mathcal{T}_{\text{cpu}}(\zeta\kappa\Pi + \Theta) \end{array} \right) \\ &= \frac{(-1)^i}{r_2 - r_1} (\mathcal{G}_{A_i} + \mathcal{G}_B\mathcal{T}_{\text{env}}(\zeta\kappa\Pi + \Theta) - \mathcal{G}_{C_i}\mathcal{T}_{\text{cpu}}(\zeta\kappa\Pi + \Theta)), \end{aligned} \quad (26)$$

where, $\mathcal{G}_{A_i} = r_{\bar{i}}\mathcal{C}_{3_{inc}}(\zeta\kappa\Pi) + \sigma_1(\mathcal{P}_{act} + k_C)$, $\mathcal{G}_B = \sigma_1 k_T$, and $\mathcal{G}_{C_i} = \beta_1 + r_{\bar{i}}$.

Similarly, from the Equation (40) in Appendix A.1, we can show that (details in the tech report),

$$\begin{aligned} &\Rightarrow \mathcal{T}_{\text{env}}(\zeta\kappa\Pi + \Theta) \left(\mathcal{J}_4(\Theta) \right) + \mathcal{T}_{\text{cpu}}(\zeta\kappa\Pi + \Theta) \left(\mathcal{J}_3(\Theta) \right) + \mathcal{T}_{\text{cpu}}(\zeta\kappa\Pi) \left(\mathcal{J}_1(\Theta) \right) \\ &+ \mathcal{T}_{\text{env}}(\zeta\kappa\Pi) \left(\mathcal{J}_2(\Theta) \right) + \mathcal{J}_A(\Theta) = 0 \end{aligned} \quad (27)$$

where, $\mathcal{J}_4(\Theta) = \mathcal{G}_B((e^{r_1(\Pi-\Theta)} - 1)(\frac{r_1+\beta_1}{k_T\sigma_1}) + (e^{r_2(\Pi-\Theta)} - 1)(\frac{r_2+\beta_1}{k_T\sigma_1}))$, $\mathcal{J}_3(\Theta) = (-\mathcal{G}_{C_1}(e^{r_1(\Pi-\Theta)} - 1)(\frac{r_1+\beta_1}{k_T\sigma_1}) - \mathcal{G}_{C_2}(e^{r_2(\Pi-\Theta)} - 1)(\frac{r_2+\beta_1}{k_T\sigma_1}))$, $\mathcal{J}_1(\Theta) = (-\mathcal{G}_{C_1}(e^{r_1\Theta} - 1)(\frac{r_1+\beta_1}{k_T\sigma_1}) - \mathcal{G}_{C_2}(e^{r_2\Theta} - 1)(\frac{r_2+\beta_1}{k_T\sigma_1}))$, $\mathcal{J}_2(\Theta) = \mathcal{G}_B((e^{r_1(\Pi-\Theta)} - 1)(\frac{r_1+\beta_1}{k_T\sigma_1}) + (e^{r_2\Theta} - 1)(\frac{r_2+\beta_1}{k_T\sigma_1}))$, and $\mathcal{J}_A(\Theta) = \mathcal{G}_{A_1}(\frac{r_1+\beta_1}{k_T\sigma_1})(e^{r_1\Theta} + e^{r_1(\Pi-\Theta)} - 2) + \mathcal{G}_{A_2}(\frac{r_2+\beta_1}{k_T\sigma_1})(e^{r_2\Theta} + e^{r_2(\Pi-\Theta)} - 2)$.

Further, we consider a CPU temperature for $(\zeta\kappa\Pi, \zeta\kappa\Pi + \Theta]$ within the stability region and find the following relationship from the Equation (37) in Appendix A.1,

$$\mathcal{T}_{\text{cpu}}^{\text{act}}(\zeta\kappa\Pi + \Theta) = \mathcal{T}_{\text{cpu}}^{\text{act}}(\zeta\kappa\Pi) + \mathcal{C}_{1_{act}}(\zeta\kappa\Pi)(e^{r_1\Theta} - 1) + \mathcal{C}_{2_{act}}(\zeta\kappa\Pi)(e^{r_2\Theta} - 1) \quad (28)$$

Substituting values for the constants from Equation (26), we get,

$$\Rightarrow \mathcal{T}_{\text{cpu}}(\zeta\kappa\Pi + \Theta) = \mathcal{T}_{\text{cpu}}^{\text{act}}(\zeta\kappa\Pi) \left(\mathcal{P}_7(\Theta) \right) + \mathcal{T}_{\text{env}}(\zeta\kappa\Pi) \left(\mathcal{P}_8(\Theta) \right) + \mathcal{P}_9(\Theta),$$

where, $\mathcal{P}_7(\Theta) = 1 - (e^{r_1\Theta} - 1)\mathcal{G}_{C_1} - (e^{r_2\Theta} - 1)\mathcal{G}_{C_2}$, $\mathcal{P}_8(\Theta) = (e^{r_2\Theta} - 1)\mathcal{G}_B + (e^{r_1\Theta} - 1)\mathcal{G}_B$, and $\mathcal{P}_9(\Theta) = (e^{r_1\Theta} - 1)\mathcal{G}_{A_1} + (e^{r_2\Theta} - 1)\mathcal{G}_{A_2}$.

Also, considering the environment thermal behavior and substituting values for the constants from Equation (26), we get (details in the tech report),

$$\Rightarrow \mathcal{T}_{\text{env}}^{\text{act}}(\zeta\kappa\Pi + \Theta) = \mathcal{T}_{\text{cpu}}(\zeta\kappa\Pi) \left(\mathcal{P}_{10}(\Theta) \right) + \mathcal{T}_{\text{env}}(\zeta\kappa\Pi) \left(\mathcal{P}_{11}(\Theta) \right) + \left(\mathcal{P}_{12}(\Theta) \right), \quad (29)$$

where, $\mathcal{P}_{10}(\Theta) = -\frac{r_1+\beta_1}{k_T\sigma_1}\mathcal{G}_{C_1}(e^{r_1\Theta} - 1) - \frac{r_2+\beta_1}{k_T\sigma_1}\mathcal{G}_{C_2}(e^{r_2\Theta} - 1)$, $\mathcal{P}_{11}(\Theta) = 1 + \frac{r_1+\beta_1}{k_T\sigma_1}\mathcal{G}_B(e^{r_1\Theta} - 1) + \frac{r_2+\beta_1}{k_T\sigma_1}\mathcal{G}_B(e^{r_2\Theta} - 1)$, and $\mathcal{P}_{12}(\Theta) = \frac{r_1+\beta_1}{k_T\sigma_1}(\mathcal{G}_{A_1}(e^{r_1\Theta} - 1) + \frac{r_2+\beta_1}{k_T\sigma_1}\mathcal{G}_{A_2}(e^{r_2\Theta} - 1))$.

Therefore, applying the Equations (25), (27), (29), and (29), in Equation (22), we may finally express our thermal-resiliency function in terms of the fixed thermal constants and input \mathcal{T}_{ref} and $\Theta^{(i)}$ (which comes from the input mode M_i) as follows,

$$\Lambda(M_i, \mathcal{T}_{\text{ref}}) = \mathcal{T}_{\text{ref}} - \frac{\mathcal{E}_1(\Theta^{(i)})}{\mathcal{E}_N(\Theta^{(i)})} - \frac{\mathcal{E}_2(\Theta^{(i)})}{\mathcal{E}_N(\Theta^{(i)})}, \quad (30)$$

where,

$$\begin{aligned} \mathcal{E}_1(\Theta) &= \mathcal{J}_A(\Theta)\mathcal{P}_2(\Theta) + \mathcal{J}_4(\Theta)\mathcal{P}_{12}(\Theta)\mathcal{P}_2(\Theta) + \mathcal{J}_A(\Theta)\mathcal{P}_{11}(\Theta)\mathcal{P}_4(\Theta) - \mathcal{J}_2(\Theta)\mathcal{P}_{12}(\Theta)\mathcal{P}_4(\Theta) \\ &+ \mathcal{J}_A(\Theta)\mathcal{P}_3(\Theta)\mathcal{P}_8(\Theta) + \mathcal{J}_4(\Theta)\mathcal{P}_{12}(\Theta)\mathcal{P}_3(\Theta)\mathcal{P}_8(\Theta) - \mathcal{J}_3(\Theta)\mathcal{P}_{12}(\Theta)\mathcal{P}_4(\Theta)\mathcal{P}_8(\Theta) \\ &- \mathcal{J}_2(\Theta)\mathcal{P}_3(\Theta)\mathcal{P}_9(\Theta) - \mathcal{J}_4(\Theta)\mathcal{P}_{11}(\Theta)\mathcal{P}_3(\Theta)\mathcal{P}_9(\Theta) + \mathcal{J}_3(\Theta)\mathcal{P}_{11}(\Theta)\mathcal{P}_4(\Theta)\mathcal{P}_9(\Theta) \\ &- \mathcal{J}_4(\Theta)\mathcal{P}_{11}(\Theta)\mathcal{P}_A(\Theta) - \mathcal{J}_3(\Theta)\mathcal{P}_8(\Theta)\mathcal{P}_A(\Theta) + \mathcal{J}_3(\Theta)\mathcal{P}_2(\Theta)\mathcal{P}_9(\Theta) - \mathcal{J}_2(\Theta)\mathcal{P}_A(\Theta), \\ \mathcal{E}_2(\Theta) &= -\mathcal{J}_A(\Theta)\mathcal{P}_1(\Theta) - \mathcal{J}_4(\Theta)\mathcal{P}_1(\Theta)\mathcal{P}_{12}(\Theta) - \mathcal{J}_A(\Theta)\mathcal{P}_{10}(\Theta)\mathcal{P}_4(\Theta) + \mathcal{J}_1(\Theta)\mathcal{P}_{12}(\Theta)\mathcal{P}_4(\Theta) \\ &- \mathcal{J}_A(\Theta)\mathcal{P}_3(\Theta)\mathcal{P}_7(\Theta) - \mathcal{J}_4(\Theta)\mathcal{P}_{12}(\Theta)\mathcal{P}_3(\Theta)\mathcal{P}_7(\Theta) + \mathcal{J}_3(\Theta)\mathcal{P}_{12}(\Theta)\mathcal{P}_4(\Theta)\mathcal{P}_7(\Theta) \\ &+ \mathcal{J}_1(\Theta)\mathcal{P}_3(\Theta)\mathcal{P}_9(\Theta) + \mathcal{J}_4(\Theta)\mathcal{P}_{10}(\Theta)\mathcal{P}_3(\Theta)\mathcal{P}_9(\Theta) - \mathcal{J}_3(\Theta)\mathcal{P}_{10}(\Theta)\mathcal{P}_4(\Theta)\mathcal{P}_9(\Theta) \\ &+ \mathcal{J}_4(\Theta)\mathcal{P}_{10}(\Theta)\mathcal{P}_A(\Theta) + \mathcal{J}_3(\Theta)\mathcal{P}_7(\Theta)\mathcal{P}_A(\Theta) - \mathcal{J}_3(\Theta)\mathcal{P}_1(\Theta)\mathcal{P}_9(\Theta) + \mathcal{J}_1(\Theta)\mathcal{P}_A(\Theta), \\ \mathcal{E}_N(\Theta) &= \mathcal{J}_2(\Theta)\mathcal{P}_1(\Theta) + \mathcal{J}_4(\Theta)\mathcal{P}_1(\Theta)\mathcal{P}_{11}(\Theta) - \mathcal{J}_1(\Theta)\mathcal{P}_2(\Theta) - \mathcal{J}_4(\Theta)\mathcal{P}_{10}(\Theta)\mathcal{P}_2(\Theta) \\ &+ \mathcal{J}_2(\Theta)\mathcal{P}_{10}(\Theta)\mathcal{P}_4(\Theta) - \mathcal{J}_1(\Theta)\mathcal{P}_{11}(\Theta)\mathcal{P}_4(\Theta) - \mathcal{J}_3(\Theta)\mathcal{P}_2(\Theta)\mathcal{P}_7(\Theta) + \mathcal{J}_2(\Theta)\mathcal{P}_3(\Theta)\mathcal{P}_7(\Theta) \\ &+ \mathcal{J}_4(\Theta)\mathcal{P}_{11}(\Theta)\mathcal{P}_3(\Theta)\mathcal{P}_7(\Theta) - \mathcal{J}_3(\Theta)\mathcal{P}_{11}(\Theta)\mathcal{P}_4(\Theta)\mathcal{P}_7(\Theta) \\ &+ \mathcal{J}_3(\Theta)\mathcal{P}_1(\Theta)\mathcal{P}_8(\Theta) - \mathcal{J}_1(\Theta)\mathcal{P}_3(\Theta)\mathcal{P}_8(\Theta) - \mathcal{J}_4(\Theta)\mathcal{P}_{10}(\Theta)\mathcal{P}_3(\Theta)\mathcal{P}_8(\Theta) \\ &+ \mathcal{J}_3(\Theta)\mathcal{P}_{10}(\Theta)\mathcal{P}_4(\Theta)\mathcal{P}_8(\Theta). \end{aligned}$$

7. VALIDATION

In this section, we evaluate our control framework both in simulations and upon an experimental hardware testbed.

7.1. Simulations

In the simulations, we simulate the execution of a single-core processor which consists of a thermal controller, PWM frequency controller loop, and scheduling algorithm. The following task parameters are used in our simulations:

- Each sporadic task $\tau_j = (e_j, d_j, p_j)$ has a period p_j uniformly drawn from the interval $[5, 15]$. (A small period range is used to keep LCM of periods from becoming too large). The execution time requirement e_j set to the task utilization times p_j , where task utilization is calculated using the UUnifast algorithm[Bini and Buttazzo 2004]. For each task, d_j equals p_j . The tasks are scheduled by EDF.
- The total number of tasks is eight; each task τ_j has three different real-time performance modes where $\tau_j^{(2)} = (e_j, d_j, p_j)$; $\tau_j^{(1)} = (.2e_j, d_j, p_j)$; and $\tau_j^{(0)}$ means that task is

Parameter	Variable	Value
CPU Active Power	\mathcal{P}_{act}	73 W
CPU Idle Power	\mathcal{P}_{inc}	20 W
Server Period	Π	20 ms
Sampling Time	T_s	100 ms
Optimal Feedback	K_o	[.5725 0]
Q matrix in Performance Index	Q	$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$
R matrix in Performance Index	R	[1]
Integral Gain	γ_I	0.00042

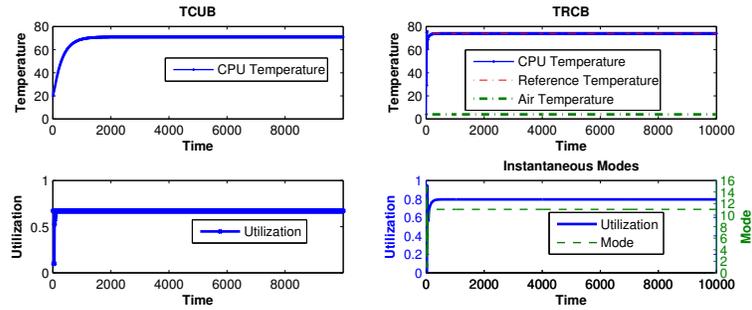


Fig. 6. Fixed \mathcal{T}_{air} for Simulation. Left plots represent TCUB and right plots represent TRCB.

not selected. From set of all possible combinations of tasks, we have selected fifteen modes with utilizations ranging from zero to one.

We refer to the controller described in Algorithm 1 as *Temperature Regulated Capacity Bound* (TRCB). In our simulations, we closely compare the performance of our proposed method with [Fu et al. 2010b] referred to as *Thermal Control Utilization Bound* (TCUB). TCUB has been chosen due to its low controller time complexity of $O(1)$. TCUB works by attempting to track a reference temperature and adjusting system utilization as needed by changing task modes via a mode assignment heuristic. The major difference between TCUB and TRCB is that TCUB does not have predefined modes. Therefore, TCUB may differ in the assigned modes from run to run for the same system temperature. Furthermore, TCUB does not use multiple power levels. TRCB on the other hand has predefined modes which permit the derivation of thermal resiliency for each mode. TRCB also utilizes a low-power mode (if available).

In our simulation, we use the same system parameters as our testbed (Intel Pentium IV 3.0 GHz). The pertinent power and control parameters are given in Table 7.1. Extensive testbed runs were carried out to generate the remaining system parameters using SI. We use the SI tools provided by Matlab to derive the system state-space parameters. Also we use the system parameters, generated from our testbed as the simulation parameters. We observe a matching of our testbed readings and the simulation. More details on this process are contained in the technical report [Hettiarachchi et al. 2011].

In Figure 6, the system response and the utilization has been shown for both TRCB (right graphs) and TCUB (left graphs) given a stable air temperature \mathcal{T}_{air} temperature equal to 5°C . The behavior of both controllers in this stable environment is nearly iden-

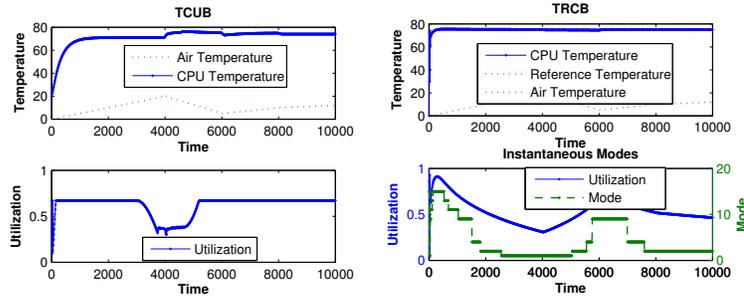


Fig. 7. Dynamically Varying T_{air} for Simulation. Left plots represent TCUB and right plots represent TRCB.

tical for thermal and utilization behavior. (The difference is due to the fact that TRCB uses EDF and TCUB uses RM scheduling). For TRCB, we also display the achieved modes at any given time in the simulation in the lower right graph.

Figure 7 shows the behavior of both TRCB and TCUB when T_{air} is dynamically changed over time. In the top two graphs of the figure, the absolute CPU temperatures over time obtained by TCUB and TRCB, respectively, are plotted along with the T_{air} . The two bottom graphs of Figure 7 present the achieved utilization for each controller; additionally, the bottom right graph displays the active mode at any point in time for TRCB. Observe that both controllers are able to track the reference temperature T_{ref} despite the sharp changes in T_{air} . For both controllers, the utilization appropriately tracks the changes in air temperature. When the air temperature increases, both controllers decrease the system utilization and increase the utilization again when the air temperature drops. Similarly, the mode plot in the lower right graph tracks the temperature changes.

Regarding the real-time performance, figures displaying deadline miss ratios have been omitted as no deadline miss was experienced for either controller in all the simulations. TCUB uses a safe utilization bound of approximately 67% to make deadline misses improbably for rate-monotonic scheduling [Liu and Layland 1973]. However, TRCB guarantees that no deadlines are ever missed due to verification using a multi-modal schedulability test [Fisher and Ahmed 2011] as described in Section 4.2.

Thus far, the empirical performance of TRCB and TCUB may appear similar. However, we believe the distinguishing feature of TRCB is the ability to guarantee hard deadlines and to calculate thermal resiliency levels during design time. Thermal resiliency calculation provides a *non-destructive* thermal stress analysis for real-time performance modes in an unpredictable operating environment. Our approach has achieved the ability to calculate the thermal resiliency by forcing the system to execute in a very predictable manner (i.e., periodic executions from PWM). To evaluate and illustrate our thermal resiliency calculation, we have used the technique in Section 6 to calculate the thermal resiliency levels for our randomly-generated multi-mode system. Figure 8 displays the thermal resiliency $\Lambda(M_i, T_{ref})$ for a range of modes and reference temperatures. Observe that the thermal resiliency increases with decreasing modes or increasing T_{ref} .

7.2. Experiments upon Hardware Testbed

To further confirm the validity of the theoretical results, we have run a task system with eight tasks, each with three modes (identical to the simulation setting), on our

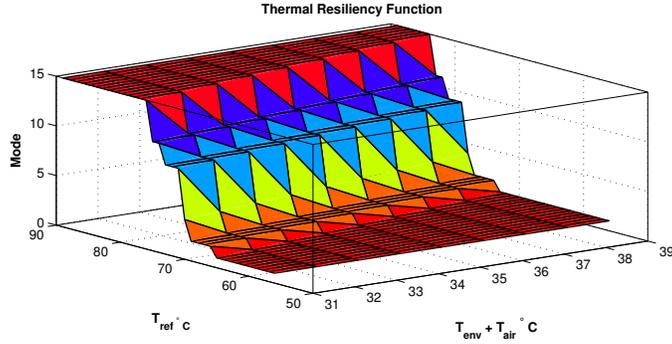


Fig. 8. Thermal resiliency over modes and \mathcal{T}_{ref} .

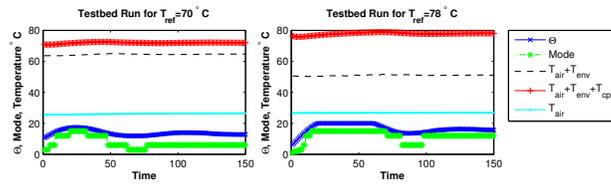


Fig. 9. The testbed running at different \mathcal{T}_{ref} values showing the Θ and Mode change over the time

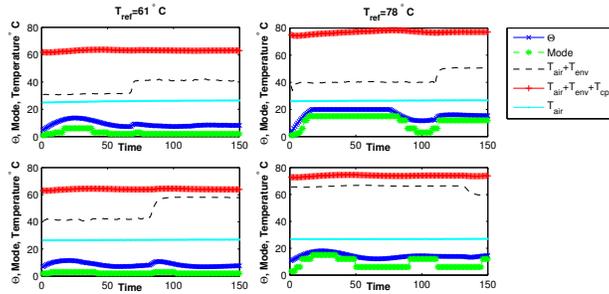


Fig. 10. The testbed running at varying environmental conditions showing the Θ and Mode change over the time

hardware testbed. Each task performs numerical calculations while executing on the system. Our hardware testbed behaves similar to the simulations of the previous subsection. Figure 9 presents testbed runs for a fixed air and environment temperature. Figure 10 shows how the testbed behaves when an outside heat source is dynamically introduced into the environment. Observe that there is a momentary drop in performance mode; however, the system soon stabilizes.

Finally, we validate our thermal resiliency calculation. Unfortunately, we do not have test equipment to accurately vary the air or environment temperature. Thus, we consider the air temperature to be fixed at the room temperature (in this case $\mathcal{T}_{air} = 24.8 \text{ } ^\circ\text{C}$). Instead, we indirectly analyze the thermal resiliency function via the inverse of the thermal resiliency function $\Lambda^{-1}(M_i, \mathcal{T}_{air}) = \min\{\mathcal{T}_{ref} \mid \mathcal{T}_{air} \leq \Lambda(M_i, \mathcal{T}_{ref})\}$. Intuitively, a lower value of $\Lambda^{-1}(M_i, \mathcal{T}_{air})$ means the system can operate at a lower temperature and thus is more resilient than a higher value of the function. We have calculated this function for four different runs of the hardware testbed (to ensure that

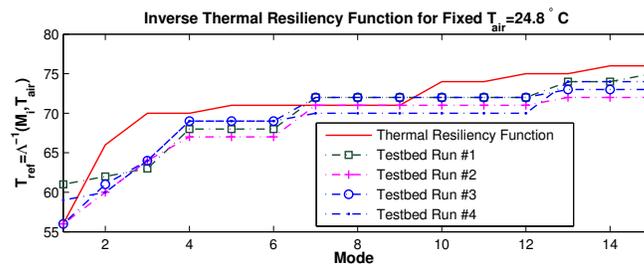


Fig. 11. Thermal Resiliency for the Simulation.

minor fluctuations of the air temperature do not affect the system). Figure 11 shows a plot of the thermal resiliency of the testbed runs when the T_{ref} is changed. The figure shows that the calculated inverse resiliency of the system increases with increasing operating mode. Most importantly, the calculated thermal resiliency tracks the actual behavior of the testbed and provides a safe upper bound on T_{ref} in a large majority of the cases which validates the effectiveness of the resiliency function.

8. CONCLUSIONS

In this paper, we have addressed the problem of obtaining performance guarantees in an unpredictable thermal environment. Towards this challenge we have presented a control-theoretic framework for thermal stress analysis in real-time systems. Our proposed method employs a nested feedback control system, which is based on optimum control theory. For our system, we derive strong thermal-resiliency and hard-real-time guarantees for any real-time performance mode. Our method has the distinct advantage of being able to verify the real-time thermal resiliency of a system before it is put into operation. In addition, we show via simulations that our framework performs as well as previous approaches which have no formal guarantee on the thermal resiliency. Our implementation upon a hardware testbed validates our proposed model and control framework.

In future work, we plan to extend our framework to control designs that are more robust to model inaccuracies (e.g., H_∞ or model-predictive controllers). As a initial step in designing a framework for thermal stress analysis, our current design uses two RC circuits (for dynamic and leakage currents) to model the CPU temperature. We plan on extending our model to permit multiple RC circuits for heterogeneous thermal distributions and generalizing our thermal equations for more complex RC circuit layouts. We hope to derive a general-theoretic design framework that captures “resiliency” metrics for other system properties (e.g., energy, noise, etc.) and extend our analysis to other hardware settings (e.g., multicore, DVS).

REFERENCES

- 2000. *Intel Pentium 4 processor in the 423-pin package thermal design guidelines*. Intel Corp.
- AHMED, M., FISHER, N., WANG, S., AND HETTIARACHCHI, P. 2011. Minimizing peak temperature in embedded real-time systems via thermal-aware periodic resources. *Sustainable Computing: Informatics and Systems 1*, 3, 226 – 240.
- BANSAL, N. AND PRUHS, K. 2005. Speed scaling to manage temperature. In *Symposium on Theoretical Aspects of Computer Science*.
- BINI, E. AND BUTTAZZO, G. 2004. Biasing effects in schedulability measures. In *Proceedings of the 16th Euromicro Conference on Real-Time Systems*. IEEE Computer Society, 196–203.
- BROOKS, D. AND MARTONOSI, M. 2001. Dynamic thermal management for high-performance microprocessors. In *International Symposium on High-Performance Computer Architecture*.

- CHANTEM, T., DICK, R. P., AND HU, X. S. 2008. Temperature-aware scheduling and assignment for hard real-time applications on MPSoCs. In *Design, Automation and Test in Europe*.
- CHEN, J.-J., HUNG, C.-M., AND KUO, T.-W. 2007. On the minimization of the instantaneous temperature for periodic real-time tasks. In *IEEE Real-Time and Embedded Technology and Applications Symposium*.
- CHEN, J.-J., WANG, S., AND THIELE, L. 2009. Proactive speed scheduling for frame-based real-time tasks under thermal constraints. In *IEEE Real-Time and Embedded Technology and Applications Symposium (RTAS)*.
- DORF, R. C. AND BISHOP, R. H. 2000. *Modern Control Systems*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA.
- FERREIRA, A., MOSSE, D., AND OH, J. 2007. Thermal faults modeling using a rc model with an application to web farms. In *Proceedings of the Euromicro Conference on Real-Time Systems*. IEEE Computer Society.
- FISHER, N. AND AHMED, M. 2011. Tractable real-time schedulability analysis for mode changes under temporal isolation. In *Proceedings of the 9th IEEE Symposium on Embedded Systems for Real-Time Multimedia (ESTImedia)*. IEEE Computer Society.
- FISHER, N., CHEN, J.-J., WANG, S., AND THIELE, L. 2009. Thermal-aware global real-time scheduling on multicore systems. In *Proceedings of the 15th IEEE Real-Time and Embedded Technology and Applications Symposium*. IEEE Computer Society Press.
- FU, X., WANG, X., AND PUSTER, E. 2010a. Simultaneous thermal and timeliness guarantees in distributed real-time embedded systems. *Journal of Systems Architecture*. To Appear.
- FU, Y., KOTTENSTETTE, N., CHEN, Y., LU, C., KOUTSOUKOS, X. D., AND WANG, H. 2010b. Feedback thermal control for real-time system. In *Proceedings of the Real-Time and Embedded Technology and Applications Systems Symposium*. IEEE Computer Society Press, Stockholm, Sweden.
- GELIG, A. K. AND CHURILOV, ALEXANDER N., . 1998. *Stability and oscillations of nonlinear pulse-modulated systems / Arkadii Kh. Gelig, Alexander N. Churilov*. Boston : Birkhauser. Includes bibliographical references (p. [343]-359) and index.
- HETTIARACHCHI, P. M., FISHER, N., AHMED, M., WANG, L. Y., WANG, S., AND SHI, W. 2011. The design and analysis of thermally-resilient hard-real-time systems (extended version). Tech. rep., Wayne State University. Available at <http://www.cs.wayne.edu/~fishern/papers/thermal-control-rtas2012.pdf>.
- KELLY, G. 2006. Body temperature variability (part 1): a review of the history of body temperature and its variability due to site selection, biological rhythms, fitness, and aging. *Alternative Medicine Review* 11, 4, 278–293.
- KIM, S., TATHIREDDY, P., NORMANN, R., AND SOLZBACHER, F. 2007. Thermal impact of an active 3-d microelectrode array implanted in the brain. *IEEE Transactions on Neural Systems and Rehabilitation Engineering* 15, 4, 493–501.
- LAMANNA, J. C., MCCRACKEN, K. A., PATIL, M., AND PROHASKA, O. J. 1989. Stimulus-activated changes in brain tissue temperature in the anesthetized rat. *Metabolic Brain Disease* 4, 4, 225–237.
- LAZZI, G. 2005. Thermal effects of bioimplants. *IEEE Engineering in Medicine and Biology Magazine* 24, 5, 75–81.
- LIU, C. AND LAYLAND, J. 1973. Scheduling algorithms for multiprogramming in a hard real-time environment. *Journal of the ACM* 20, 1, 46–61.
- LIU, Y., DICK, R. P., SHANG, L., AND YANG, H. 2007. Accurate temperature-dependent integrated circuit leakage power estimation is easy. In *Proceedings of the conference on Design, automation and test in Europe*. Nice, France, 1526–1531.
- MOK, A. K. 1983. Fundamental design problems of distributed systems for the hard-real-time environment. Ph.D. thesis, Laboratory for Computer Science, Massachusetts Institute of Technology. Available as Technical Report No. MIT/LCS/TR-297.
- NISE, N. S. 2000. *Control Systems Engineering*. John Wiley & Sons, Inc., New York, NY, USA.
- OGATA, K. 1995. *Discrete-time control systems (2nd ed.)*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA.
- QUAN, G. AND ZHANG, Y. 2009. Leakage Aware Feasibility Analysis for Temperature-Constrained Hard Real-Time Periodic Tasks. In *Proceedings of the 2009 21st Euromicro Conference on Real-Time Systems-Volume 00*. IEEE Computer Society, 207–216.

- RAJKUMAR, R., LEE, C., LEHOCZKY, J., AND SIEWIOREK, D. 1997. A resource allocation model for qos management. In *Proceedings of the 18th IEEE Real-Time Systems Symposium*. RTSS '97. IEEE Computer Society, Washington, DC, USA, 298–.
- RUGGERA, P., WITTERS, D., VON MALTZAHN, G., AND BASSEN, H. 2003. *In vitro* assessment of tissue heating near metallic medical implants by exposure to pulsed radio frequency diathermy. *Physics in Medicine and Biology* 48, 17, 2919–2928.
- SERGENT, J. AND KRUM, A. 1998. *Thermal Management Handbook for Electronic Assemblies*. McGraw-Hill Professional.
- SHIN, I. AND LEE, I. 2008. Compositional real-time scheduling framework with periodic model. *ACM Transactions on Embedded Computing Systems* 7, 3.
- SKADRON, K., STAN, M. R., HUANG, W., VELUSAMY, S., SANKARANARAYANAN, K., AND TARIJAN, D. 2003. Temperature-aware microarchitecture. In *International Symposium on Computer Architecture*.
- TIMMONS, N. AND SCANLON, W. 2009. An adaptive energy efficient mac protocol for the medical body area network. In *1st International Conference on Wireless Communication, Vehicular Technology, Information Theory and Aerospace Electronic Systems Technology, 2009*. 587 – 593.
- WANG, S. AND BETTATI, R. 2008. Reactive speed control in temperature-constrained real-time systems. *Real-Time Systems Journal* 39, 1-3, 658–671.

A. APPENDIX

A.1. The Temperature Calculations

When we consider our thermal model with the leakage current effect, the CPU temperature is calculated based on the solution of second order differential equation. From Equation (3), we get the first derivative of $\mathcal{T}_{\text{env}}(t)$ as follows,

$$\begin{aligned} \frac{d}{dt}\mathcal{T}_{\text{env}}(t) &= \frac{1}{k_T\sigma_1} \left(\frac{d^2}{dt^2}\mathcal{T}_{\text{cpu}}(t) + \beta_1 \frac{d}{dt}\mathcal{T}_{\text{cpu}}(t) - \sigma_1 \frac{d}{dt}\mathcal{P}_{\text{cpu}}^d(t) \right) \\ &= \frac{1}{k_T\sigma_1} \left(\frac{d^2}{dt^2}\mathcal{T}_{\text{cpu}}(t) + \beta_1 \frac{d}{dt}\mathcal{T}_{\text{cpu}}(t) \right). \end{aligned} \quad (31)$$

In this analysis, we consider a system that can be described according to the model shown in the Section 4. Therefore, in the above Equation (31), we consider the system behavior for discrete time intervals and the input is considered to be constant in each sampling interval (the input value at the sampling time continue to hold for the rest of the period, until the next sampling time). This assumption is realistic because we implement our system as a discrete-time control system, in which the ZOH functionality means for holding the input value for inter-sampling times periods. Let us consider any such general time period where the input is held constant; therefore, for time instant t in this range, $\frac{d}{dt}\mathcal{P}_{\text{cpu}}^d(t)$ can be considered as zero. Thus, we can substitute Equation (2) and (31) to the Equation (5) to get the following,

$$\frac{d^2}{dt^2}\mathcal{T}_{\text{cpu}}(t) + \mathcal{V} \frac{d}{dt}\mathcal{T}_{\text{cpu}}(t) + \mathcal{B}\mathcal{T}_{\text{cpu}}(t) = \mathcal{F}_{\text{act/inc/cont}}, \quad (32)$$

where,

$$\begin{aligned} \mathcal{V} &\stackrel{\text{def}}{=} (\beta_1 + \beta_2), \\ \mathcal{B} &\stackrel{\text{def}}{=} (\beta_1\beta_2 - k_T^2\sigma_1\sigma_2), \\ \mathcal{F}_{\text{act/inc/cont}} &\stackrel{\text{def}}{=} (\beta_2\sigma_1 + \sigma_1\sigma_2k_T) (\mathcal{P}_{\text{act/inc/cont}} + k_C) + \sigma_1\sigma_2k_T\mathcal{P}_{\text{env}}(t), \\ \mathcal{P}_{\text{act/inc/cont}} &= \begin{cases} \mathcal{P}_{\text{act}}, & \text{active;} \\ \mathcal{P}_{\text{inc}}, & \text{inactive;} \\ (\mathcal{P}_{\text{act}} - \mathcal{P}_{\text{inc}}) \frac{\Theta}{\Pi} + \mathcal{P}_{\text{inc}}, & \text{continuous.} \end{cases} \end{aligned} \quad (33)$$

The Equation (32) is a second-order inhomogeneous equation and \mathcal{F} is a constant ($\mathcal{P}_{\text{cpu}}^{\text{d}}(t)$ and $\mathcal{P}_{\text{env}}(t)$ are unchanged over two sampling periods). As we already discuss, the CPU can operate in two power modes. Depending on the operating mode of the system (active or inactive CPU operation), we can derive two different \mathcal{F} values. Also, we assume, when the CPU power is represented in terms of the resource capacity, Θ , the corresponding \mathcal{F} is denoted by $\mathcal{F}_{\text{cont}}$.⁵ Therefore, the complete solution for \mathcal{T}_{cpu} and \mathcal{T}_{env} over any continuous interval is given by,

$$\mathcal{T}_{\text{cpu}}^{\text{act/inc}}(t) = \mathcal{C}_{1_{\text{act/inc}}} e^{r_1 t} + \mathcal{C}_{2_{\text{act/inc}}} e^{r_2 t} + \mathcal{C}_{3_{\text{act/inc}}}, \quad (34)$$

where, $r_{1/2} = -\frac{1}{2}(\mathcal{V} \mp \sqrt{\mathcal{V}^2 - 4\mathcal{B}})$ and $\mathcal{C}_{3_{\text{act/inc/cont}}} = \frac{\mathcal{F}_{\text{act/inc/cont}}}{\mathcal{B}}$.

In the Equation (34), the r_1 and r_2 terms are negative because $\sqrt{\mathcal{V}^2 - 4\mathcal{B}}$ is positive and less than \mathcal{V} .

From Equation (5) and (34), we can find the $\mathcal{T}_{\text{env}}(t)$ for active and inactive CPU operations as follows,

$$\begin{aligned} \mathcal{T}_{\text{env}}^{\text{act/inc}}(t) &= \frac{1}{k_T \sigma_1} \left(\mathcal{C}_{1_{\text{act/inc}}} r_1 e^{r_1 t} + \mathcal{C}_{2_{\text{act/inc}}} r_2 e^{r_2 t} - \sigma_1 (\mathcal{P}_{\text{act/inc}} + k_c) \right) \\ &+ \frac{\beta_1}{k_T \sigma_1} \left(\mathcal{C}_{1_{\text{act/inc}}} e^{r_1 t} + \mathcal{C}_{2_{\text{act/inc}}} e^{r_2 t} + \mathcal{C}_{3_{\text{act/inc}}} \right). \end{aligned} \quad (35)$$

We consider the system operates in interleaved active and inactive power modes over given interval size; the initial temperature of given period is the final temperature of the previous period. Given $\mathcal{T}_{\text{cpu}}(t_b)$ and $\mathcal{T}_{\text{env}}(t_b)$, fixed $\mathcal{P}_{\text{cpu}}^{\text{d}}$ and \mathcal{P}_{env} , we may obtain $\mathcal{C}_1, \mathcal{C}_2$ by solving Equations (34) and (35), where t_b is the initial time of the interval. Further, we derive the \mathcal{C} as follows,

$$\begin{aligned} \mathcal{C}_{1_{\text{act/inc/cont}}}(t_b) &= \frac{1}{r_1 - r_2} \left(r_2 \mathcal{C}_{3_{\text{act/inc}}} + \sigma_1 (\mathcal{P}_{\text{act/inc/cont}} + k_C + k_T \mathcal{T}_{\text{env}}(t_b)) - (\beta_1 + r_2) \mathcal{T}_{\text{cpu}}(t_b) \right), \\ \mathcal{C}_{2_{\text{act/inc/cont}}}(t_b) &= \frac{1}{r_2 - r_1} \left(r_1 \mathcal{C}_{3_{\text{act/inc}}} + \sigma_1 (\mathcal{P}_{\text{act/inc/cont}} + k_C + k_T \mathcal{T}_{\text{env}}(t_b)) - (\beta_1 + r_1) \mathcal{T}_{\text{cpu}}(t_b) \right), \\ \mathcal{C}_{3_{\text{act/inc/cont}}}(t_b) &= \frac{(\beta_2 \sigma_1 + \sigma_1 \sigma_2 k_T) (\mathcal{P}_{\text{act/inc/cont}} + k_C) + \sigma_1 \sigma_2 k_T \mathcal{P}_{\text{env}}(t_b)}{\beta_1 \beta_2 - k_T^2 \sigma_1 \sigma_2}. \end{aligned} \quad (36)$$

Note that, here we replace the initial power settings $\mathcal{P}_{\text{cpu}}^{\text{d}}(t)$ with \mathcal{P}_{act} . We use the Equation (34) and (35) to derive the temperature of the system at the end of each period. Therefore, consider the CPU temperature at any active period, $(n\Pi, n\Pi + \Theta]$ and adjacent inactive $(n\Pi + \Theta, (n+1)\Pi]$ period (details given in the tech report).

$$\mathcal{T}_{\text{cpu}}^{\text{act}}(n\Pi + \Theta) = \mathcal{T}_{\text{cpu}}^{\text{act}}(n\Pi) + \mathcal{C}_{1_{\text{act}}}(n\Pi)(e^{r_1 \Theta} - 1) + \mathcal{C}_{2_{\text{act}}}(n\Pi)(e^{r_2 \Theta} - 1) \quad (37)$$

$$\begin{aligned} \mathcal{T}_{\text{cpu}}^{\text{inc}}((n+1)\Pi) &= \mathcal{T}_{\text{cpu}}^{\text{act}}(n\Pi) + \mathcal{C}_{1_{\text{inc}}}(n\Pi + \Theta)(e^{r_1(\Pi - \Theta)} - 1) + \mathcal{C}_{2_{\text{inc}}}(n\Pi + \Theta)(e^{r_2(\Pi - \Theta)} - 1) \\ &+ \mathcal{C}_{1_{\text{act}}}(n\Pi)(e^{r_1 \Theta} - 1) + \mathcal{C}_{2_{\text{act}}}(n\Pi)(e^{r_2 \Theta} - 1) \end{aligned} \quad (38)$$

⁵We justify the need for $\mathcal{F}_{\text{cont}}$ in the Tech report) under PWM Error Calculation

Therefore, we can derive the equation for the period $(n\Pi, (n + \varsigma)\Pi]$ is as follows.

$$\begin{aligned}
\mathcal{T}_{\text{cpu}}^{\text{inc}}((n + \varsigma)\Pi) &= \sum_{i=0}^{\varsigma-1} \mathcal{C}_{1_{\text{inc}}}((n + i)\Pi + \Theta)(e^{r_1(\Pi-\Theta)} - 1) + \sum_{i=0}^{\varsigma-1} \mathcal{C}_{2_{\text{inc}}}((n + i)\Pi + \Theta)(e^{r_2(\Pi-\Theta)} - 1) \\
&+ \mathcal{T}_{\text{cpu}}^{\text{act}}(n\Pi) + \sum_{i=0}^{\varsigma-1} \mathcal{C}_{1_{\text{act}}}((n + i)\Pi)(e^{r_1\Theta} - 1) + \sum_{i=0}^{\varsigma-1} \mathcal{C}_{2_{\text{act}}}((n + i)\Pi)(e^{r_2\Theta} - 1), \\
&= \mathcal{T}_{\text{cpu}}^{\text{act}}(n\Pi) \\
&+ \sum_{i=0}^{\varsigma-1} \sum_{j=1}^2 \left(\mathcal{C}_{(j)_{\text{inc}}}((n + i)\Pi + \Theta)(e^{r^{(j)}(\Pi-\Theta)} - 1) + \mathcal{C}_{(j)_{\text{act}}}((n + i)\Pi)(e^{r^{(j)}\Theta} - 1) \right) \quad (39)
\end{aligned}$$

Please note that the above Equation (39) is inductively defined, as the constants for the boundary conditions can be derived from the Equation (36) which are in terms of previous values of \mathcal{T}_{cpu} and \mathcal{T}_{env} .

Now we use the same approach to derive the environment temperature \mathcal{T}_{env} (the details given in the tech report) and derive the following for the period $(n\Pi, n\Pi + \varsigma]$ is as follows,

$$\begin{aligned}
\mathcal{T}_{\text{env}}^{\text{inc}}((n + \varsigma)\Pi) &= \mathcal{T}_{\text{env}}^{\text{act}}(n\Pi) + \sum_{i=0}^{\varsigma-1} \mathcal{C}_{1_{\text{inc}}}((n + i)\Pi + \Theta)(e^{r_1(\Pi-\Theta)} - 1) \frac{r_1 + \beta_1}{k_T \sigma_1} \\
&\sum_{i=0}^{\varsigma-1} \mathcal{C}_{2_{\text{inc}}}((n + i)\Pi + \Theta)(e^{r_2(\Pi-\Theta)} - 1) \frac{r_2 + \beta_1}{k_T \sigma_1} + \sum_{i=0}^{\varsigma-1} \mathcal{C}_{1_{\text{act}}}((n + i)\Pi)(e^{r_1\Theta} - 1) \frac{r_1 + \beta_1}{k_T \sigma_1} \\
&+ \sum_{i=0}^{\varsigma-1} \mathcal{C}_{2_{\text{act}}}((n + i)\Pi)(e^{r_2\Theta} - 1) \frac{r_2 + \beta_1}{k_T \sigma_1} \\
&= \mathcal{T}_{\text{env}}^{\text{act}}(n\Pi) + \sum_{i=0}^{\varsigma-1} \sum_{j=1}^2 \left(\mathcal{C}_{(j)_{\text{inc}}}((n + i)\Pi + \Theta)(e^{r^{(j)}(\Pi-\Theta)} - 1) \frac{r^{(j)} + \beta_1}{k_T \sigma_1} \right. \\
&\left. + \mathcal{C}_{(j)_{\text{act}}}((n + i)\Pi)(e^{r^{(j)}\Theta} - 1) \frac{r^{(j)} + \beta_1}{k_T \sigma_1} \right). \quad (40)
\end{aligned}$$

In the Equation (40), the constants for the boundary conditions can be derived from the Equation (36).

The CPU and environment temperature calculation equations (Equation (39) and (40)) gives a possible way to calculate the temperature states of the system, provided that we know a single boundary condition. Therefore, when we calculate the thermal resiliency and the PWM error for second-order thermal model, Equation (39) and (40) are used.

A.2. Calculation of State-Space Parameters Using Testbed Results

The state-space parameter generation process needs input and output data collected over sufficiently larger period. Unlike the testbench output ($\mathcal{T}_{\text{cpu}} + \mathcal{T}_{\text{env}}$ reading-measured using T-type thermocouple as explained in the Section 4.1), the testbed input, the equivalent CPU thermal input power cannot be measured directly. Instead, we measure the CPU input power, the closest measurable parameter. We assume the electrical power consumed by the CPU totally converts to thermal energy and measure the CPU input power and consider it as the equivalent thermal power⁶.

⁶This assumption is realistic because in the CPU (any electrical circuit) the desired objective is to operate its switches. However each gate (in switches) consume energy and generate heat. There is no any other energy transformation in an ordinary electrical circuit.

We install two shunt resistors in series with the 4-pin ATX power connector and measure the voltage (and calculate the current drawn) drop across it using National Instrument data acquisition interface, NI 9205. Since The NI 9205 does not have a Linux USB driver, we create an application interface in a Windows computer to connect with the testbed using the Ethernet. The testbed measures the CPU and environment temperature and sends a sync signal to Windows computer with NI 9205 interface to record the ATX current readings. We calculate the the total power fed to the CPU, as the current drawn by CPU (through the NI measurements) and the voltage of the 4 wire ATX interface are known.

We run a random workload for a longer time period to generate thermal effects on the CPU and record input (power) and output (CPU temperature) data. We collect two sets of data from the testbed, one set to generate the model parameters and the other set validates them. We use standard tools provided by system-identification toolbox in Matlab to derive the state-space parameters (SSP) with the test data.⁷ We use these SSP in the rest of the simulations and in the controller design.

We observe that when we do the SI process, the thermal output of the CPU is not sufficient enough to make a accurately measurable temperature difference in the environment. Therefore, for the parameter generation purpose, we consider the following: we use a first order CPU thermal-model, for the parameter generation, considering that the system environmental temperature stays stable and the the thermal model of the system is considered as a differential model. In other words, the leakage power of the testbed is a constant for a given temperature and, therefore when we consider the differential model (the difference between any steady point to the current point), the leakage power component need not to be considered for closer operational points.

When we consider the environment temperature is nearly stable over a sufficiently larger time period, we may get a normalized thermal model of the CPU as follows,

$$\frac{d}{dt}\mathcal{T}_{\text{cpu}}(t) = \sigma_1 \left(k_T - \frac{1}{R_{\text{cpu}}^l} - \frac{1}{R_{\text{cpu}}^d} \right) \mathcal{T}_{\text{cpu}}(t) + k_T \sigma_1 \mathcal{T}_{\text{env}}(t) + \sigma_1 \mathcal{P}_{\text{cpu}}^d(t) + \sigma_1 k_C. \quad (41)$$

Consider an another test-point (at t_E) during our SI process and assume the same environmental temperature, then we get, the following differential system model,

$$\frac{d}{dt}\bar{\mathcal{T}}_{\text{cpu}}(t) = A\bar{\mathcal{T}}_{\text{cpu}}(t) + B\bar{\mathcal{P}}_{\text{cpu}}^d(t), \quad (42)$$

where, $\bar{\mathcal{T}}_{\text{cpu}}(t) = \mathcal{T}_{\text{cpu}}(t) - \mathcal{T}_{\text{cpu}}(t_E)$, and $\bar{\mathcal{P}}_{\text{cpu}}^d(t) = \mathcal{P}_{\text{cpu}}^d(t) - \mathcal{P}_{\text{cpu}}^d(t_E)$. Therefore, the final system that we used in the controller design and the parameter generation may be considered as the above model.

In our parameter generation process, we use the discrete form of the above state-space Equation (42). As we shown earlier, the continuous-time state-space model can converted to discrete-time state-space model and the following discrete model is obtained,

$$\bar{\mathcal{T}}_{\text{cpu}}(k+1) = G\bar{\mathcal{T}}_{\text{cpu}}(k) + H\bar{\mathcal{P}}_{\text{cpu}}^d(k). \quad (43)$$

This parameter generation can be considered as linearization of our model at the operating points (at a particular environment temperature point). In our future work, we will generate linearized system parameters for a smooth operating regions and will implement a gain scheduled controller.

⁷We use Predictive Error Method (PEM) algorithm implementation in Matlab.