

# Power-Aware CPU Utilization Control for Distributed Real-Time Systems\*

Xiaorui Wang<sup>†</sup>, Xing Fu<sup>†</sup>, Xue Liu<sup>‡</sup>, and Zonghua Gu<sup>§</sup>

<sup>†</sup>Department of Electrical Engineering and Computer Science, University of Tennessee, Knoxville, USA,

<sup>‡</sup>McGill University, Canada, and <sup>§</sup> Hong Kong University of Science and Technology.

{xwang, xful}@utk.edu xueliu@cs.mcgill.ca zgu@cse.ust.hk

## Abstract

CPU utilization control has recently been demonstrated to be an effective way of meeting end-to-end deadlines for distributed real-time systems running in unpredictable environments. However, current research on utilization control focuses exclusively on task rate adaptation, which cannot effectively handle rate saturation and discrete task rates. Since the CPU utilization contributed by a real-time periodic task is determined by both its rate and execution time, CPU frequency scaling can be used to adapt task execution times for power-efficient utilization control. In this paper, we present a two-layer coordinated CPU utilization control architecture. The primary control loop uses frequency scaling to locally control the CPU utilization of each processor, while the secondary control loop adopts rate adaptation to control the utilizations of all the processors at the cluster level on a finer timescale. Both the two control loops are designed and coordinated based on well-established control theory for theoretically guaranteed control accuracy and system stability. Empirical results on a physical testbed demonstrate that our control solution outperforms a state-of-the-art utilization control algorithm by having more accurate control and less power consumption.

## 1 Introduction

Traditional approaches to handling end-to-end real-time tasks, such as end-to-end scheduling [1] and distributed priority ceiling [2], rely on schedulability analysis, which requires *a priori* knowledge of the tasks' Worst-Case Execution Times (WCET). While such open-loop approaches work effectively in the closed execution environments of traditional real-time systems, they may violate the desired timing constraints or severely underutilize the system when task execution times are highly unpredictable. In recent years, a new category of real-time applications called Distributed Real-time Embedded (DRE) systems has been rapidly growing. DRE systems commonly execute in open and *unpredictable* environments in which workloads are unknown and may vary significantly at runtime. Such systems include data-driven systems whose execution is heavily in-

fluenced by volatile environments. For example, task execution times in vision-based feedback control systems depend on the content of live camera images of changing environments [3]. DRE systems call for a paradigm shift from classical real-time computing that relies on accurate characterization of workloads and platform.

Recently, feedback control techniques have shown a lot of promise in providing real-time guarantees for DRE systems by adapting to workload variations based on dynamic feedback. In particular, feedback-based CPU utilization control [4][5] has been demonstrated to be an effective way of meeting the *end-to-end deadlines* for soft DRE systems. The primary goal of utilization control is to enforce appropriate schedulable utilization bounds (e.g., the Liu and Layland bound for RMS) on all the processors in a DRE system, despite significant uncertainties in system workloads. In the meantime, it tries to maximize the system utility by controlling CPU utilizations to stay slightly below their schedulable bounds so that the processors can be utilized to the maximum degree. Utilization control can also enhance system survivability by providing overload protection against workload fluctuation [6].

However, previous research on CPU utilization control exclusively relies on task rate adaptation by assuming that task rates can be continuously tuned within specified ranges. While rate adaptation is an effective actuator for some DRE systems, it has several limitations. First, it is often infeasible to achieve desired utilization set points by rate adaptation alone [7]. For example, many DRE systems are configured based on tasks' WCETs. Consequently, even when all the tasks are running at their highest rates, CPU utilizations are still way below the desired set points, resulting in severely underutilized systems and excessive power consumption. In that case, CPU frequency scaling can be used for power savings while keeping the utilizations slightly below the schedulable bounds. Second, many tasks in DRE systems only support a few discrete rates. While optimization strategies [8][9] are developed to handle discrete task rates, they rely on the common assumption that task WCETs are known *a priori* and accurate, which makes them less applicable to DRE systems running in *unpredictable* environments. Third, the model of task rate in many applications could be complex and vary at runtime based on application evolution [10][11]. As a result, the estimated task rate ranges are often inaccurate and may change sig-

\*This work was supported, in part, by NSF CSR Grant CNS-0720663, NSF CAREER Award CNS-0845390, and a Power-Aware Computing Award from Microsoft Research in 2008.

nificantly online, leading to unexpected rate saturation and even deadline misses when utilizations are higher than the schedulable bounds and can be lowered down only by rate adaptation. Finally, some DRE systems may not allow rate adaptation for any tasks but their CPU utilizations still need to be controlled. Therefore, it is important to explore complementary ways for effective CPU utilization control.

In this paper, we propose to use Dynamic Voltage and Frequency Scaling (DVFS) jointly with rate adaptation for utilization control. Since the CPU utilization contributed by a real-time periodic task is determined by both its rate and its execution time, CPU frequency scaling can be used to adapt task execution time for power-efficient utilization control. The integration of DVFS in utilization control introduces several new challenges. First, a centralized controller for simultaneous rate adaptation and DVFS would have a Multi-Input-Multi-Output (MIMO) nonlinear model. Therefore, multiple linear control loops are more preferable for acceptable runtime overhead. Second, different control loops need to be carefully designed to coordinate together for the desired control functions. Finally, the control accuracy and global system stability of the coordinated control solution must be analytically assured.

This paper presents a two-layer coordinated CPU utilization control architecture. The primary control loop uses DVFS to locally control the CPU utilization of each processor. In the meantime, the secondary control loop adopts rate adaptation to control the utilizations of all the processors at the cluster level on a finer timescale. Specifically, the contributions of this paper are four-fold:

- We derive an analytical model that captures the system dynamics of the new CPU utilization control problem.
- We design a two-layer coordinated control architecture and conduct detailed coordination analysis.
- We implement our control architecture in an open-source real-time middleware system.
- We present empirical results to demonstrate that our control solution outperforms a state-of-the-art utilization controller that relies solely on rate adaptation.

The rest of this paper is organized as follows. We formulate the new CPU utilization control problem in Section 2. Section 3 presents the system model and control architecture. Section 4 briefly introduces the rate adaptation loop while Section 5 provides the detailed design and analysis of the CPU frequency scaling loop. Section 6 discusses the implementation of the control architecture in a real-time middleware system. Section 7 presents our empirical results on a physical testbed. Section 8 reviews the related work. Finally, Section 9 summarizes the paper.

## 2 Problem Formulation

In this section, we formulate the new CPU utilization control problem for DRE systems.

### 2.1 Task Model

We adopt an end-to-end task model [12] implemented by many DRE applications. A system is comprised of  $m$  periodic tasks  $\{T_i | 1 \leq i \leq m\}$  executing on  $n$  processors  $\{P_i | 1 \leq i \leq n\}$ . Task  $T_i$  is composed of a set of subtasks  $\{T_{ij} | 1 \leq j \leq n_i\}$  which may be located on different processors. A processor may host one or more subtasks of a task. The release of subtasks is subject to precedence constraints, i.e., subtask  $T_{ij}$  ( $1 < j \leq n_i$ ) cannot be released for execution until its predecessor subtask  $T_{ij-1}$  is completed. All the subtasks of a task share the same rate. The rate of a task (and all its subtasks) can be adjusted by changing the rate of its first subtask. If a non-greedy synchronization protocol (e.g., release guard [1]) is used to enforce the precedence constraints, every subtask are released periodically without jitter.

In our task model, each task  $T_i$  has a *soft* end-to-end deadline related to its period. In an end-to-end scheduling approach [1], the deadline of an end-to-end task is divided into subdeadlines of its subtasks. Hence the problem of meeting the end-to-end deadline can be transformed to the problem of meeting the subdeadline of each subtask. A well known approach for meeting the subdeadlines on a processor is to ensure its utilization remains below its schedulable utilization bound [12].

Our task model has three important properties. First, while each subtask  $T_{ij}$  has an *estimated* execution time  $c_{ij}$  available at design time, its *actual* execution time may be different from its estimation and vary at run-time due to two reasons: CPU frequency scaling or workload uncertainties. Modeling such uncertainties is important to DRE systems operating in unpredictable environments. Second, the rate of a task  $T_i$  may be dynamically adjusted within a range  $[R_{min,i}, R_{max,i}]$ . This assumption is based on the fact that the task rates in many applications (e.g., digital control [13], sensor update, and multimedia [14]) can be dynamically adjusted without causing system failure. The rate ranges are determined by the applications (e.g., the limited sampling frequency of a sensor) and are not necessarily accurate. A task running at a higher rate contributes a higher value to the application at the cost of higher utilizations. Please note that our solution does *not* rely on continuous task rates. For a task with only discrete rates, its continuous rate value will be truncated to the highest discrete rate supported by the task that is below the continuous value. The utilization difference resulted from the truncation can be compensated by CPU frequency scaling. Third, the CPU frequency of each processor  $P_i$  may be dynamically adjusted within a range  $[F_{min,i}, F_{max,i}]$ . This assumption is based on the fact that many today's processors are DVFS-enabled. For processors that do not support DVFS, clock modulation can be used instead to change CPU frequency [15]. The frequency ranges are assumed to be continuous because a continuous value can be approximated by a series of discrete frequency levels supported by a processor, as we explain in Section 6.

## 2.2 Problem Formulation

Utilization control can be formulated as a dynamic constrained optimization problem. We first introduce some notation.  $T_s$ , the control period, is selected so that multiple instances of each task may be released during a control period.  $u_i(k)$  is the CPU utilization of processor  $P_i$  in the  $k^{th}$  control period, i.e., the fraction of time that  $P_i$  is not idle during time interval  $[(k-1)T_s, kT_s)$ .  $B_i$  is the desired utilization set point on  $P_i$ .  $r_j(k)$  is the invocation rate of task  $T_j$  in the  $(k+1)^{th}$  control period.  $f_i(k)$  is the relative CPU frequency (i.e., CPU frequency relative to the highest level  $F_{max,i}$ ) of processor  $P_i$  in the  $(k+1)^{th}$  control period.

Given a utilization set-point vector,  $\mathbf{B} = [B_1 \dots B_n]^T$ , rate constraints  $[R_{min,j}, R_{max,j}]$  for each task  $T_j$ , and frequency constraints  $[F_{min,i}, F_{max,i}]$  for each processor  $P_i$ , the control goal at  $k^{th}$  sampling point (time  $kT_s$ ) is to dynamically choose task rates  $\{r_j(k) | 1 \leq j \leq m\}$  and CPU frequencies  $\{f_i(k) | 1 \leq i \leq n\}$  to minimize the difference between  $B_i$  and  $u_i(k)$  for all the processors:

$$\min_{\{r_j(k) | 1 \leq j \leq m, f_i(k) | 1 \leq i \leq n\}} \sum_{i=1}^n (B_i - u_i(k+1))^2 \quad (1)$$

subject to constraints

$$R_{min,j} \leq r_j(k) \leq R_{max,j} \quad (1 \leq j \leq m) \quad (2)$$

$$F_{min,i} \leq f_i(k) \leq F_{max,i} \quad (1 \leq i \leq n) \quad (3)$$

The rate constraints ensure all tasks remain within their acceptable rate ranges. The frequency constraints ensure all CPU frequencies remain within their acceptable ranges. The optimization formulation minimizes the difference between the utilization of each processor and its corresponding set point, by manipulating the rate of every task and the frequency of every processor within their constraints. The design goal is to ensure that all processors quickly converge to their utilization set points after a workload variation, whenever it is feasible under the constraints. Therefore, to guarantee end-to-end deadlines, a user only needs to specify the set point of each processor to be a value below its schedulable utilization bound. Utilization control algorithms can be used to meet all the end-to-end deadlines by enforcing the set points of all the processors in a DRE system, when feasible under the constraints<sup>1</sup>.

## 3 End-to-End Utilization Control

In this section, we model the end-to-end utilization control problem and present our two-layer control architecture.

### 3.1 System Modeling

Following a control-theoretic methodology, we establish a dynamic model that characterizes the relationship between

<sup>1</sup>A system must apply admission control when its load exceeds the limit that can be handled within the rate and frequency constraints.

the controlled variable  $\mathbf{u}(k)$  and the manipulated variables  $\mathbf{r}(k)$  and  $\mathbf{f}(k)$ . We first model the utilization  $u_i(k)$  of one processor  $P_i$ . As observed in previous research [16][17], the execution times of tasks on  $P_i$  can be approximately estimated to be a linear function of  $P_i$ 's relative CPU frequency<sup>2</sup>. Therefore, the *estimated* execution time of task  $T_{jl}$  in the  $k^{th}$  control period can be modeled as  $c_{jl}/f_i(k)$ . The *estimated* CPU utilization of processor  $P_i$  can be modeled as:

$$b_i(k) = \frac{\sum_{T_{jl} \in S_i} c_{jl} r_j(k)}{f_i(k)} \quad (4)$$

where  $S_i$  is the set of subtasks located at processor  $P_i$ .

**Example:** Consider a system with two processors and three tasks.  $T_1$  has only one subtask  $T_{11}$  on processor  $P_1$ .  $T_2$  has two subtasks  $T_{21}$  and  $T_{22}$  on processors  $P_1$  and  $P_2$ , respectively.  $T_3$  has one subtask  $T_{31}$  allocated to processors  $P_2$ . The estimated utilizations of  $P_1$  and  $P_2$  are:

$$b_1(k) = \frac{c_{11} r_1(k) + c_{21} r_2(k)}{f_1(k)}$$

$$b_2(k) = \frac{c_{22} r_2(k) + c_{31} r_3(k)}{f_2(k)}$$

Note that the utilizations of  $P_1$  and  $P_2$  are coupled because the task rate of  $T_2$ , i.e.,  $r_2(k)$ , affects the utilizations of both  $P_1$  and  $P_2$ . We then define the *estimated* utilization change of  $P_i$ ,  $\Delta b_i(k)$ , as:

$$\Delta b_i(k) = \frac{\sum_{T_{jl} \in S_i} c_{jl} r_j(k)}{f_i(k)} - \frac{\sum_{T_{jl} \in S_i} c_{jl} r_j(k-1)}{f_i(k-1)} \quad (5)$$

Note that  $\Delta b_i(k)$  is based on the estimated execution time  $c_{jl}$ . Since the actual execution times may be different from their estimation due to workload variations, we model the actual utilization of  $P_i$ ,  $u_i(k)$ , as the following difference equation.

$$u_i(k+1) = u_i(k) + g_i \Delta b_i(k) \quad (6)$$

where the utilization gain  $g_i$  represents the ratio between the change to the actual utilization and its estimation  $\Delta b_i(k)$ . For example,  $g_i = 2$  means that the actual change to utilization is twice the estimated change. Note that the exact value of  $g_i$  is *unknown* at design time due to the unpredictability of subtasks' execution times.

The system model (6) is nonlinear because of the definition of  $\Delta b_i(k)$  in (5). Therefore, we need linearization to simplify the controller design for acceptable runtime overhead. There are two ways to linearize the system model. First, we may assume that all the processors always run at their highest CPU frequency and the utilizations are controlled by rate adaptation only. As a result,  $f_i(k)$  becomes 1 and the system model (6) becomes a linear model between  $\Delta b_i(k)$  and  $\Delta r_j(k) = r_j(k) - r_j(k-1)$ . Second, we can assume that the utilizations are controlled by

<sup>2</sup>In general, the execution times of some tasks may include frequency-independent parts that do not scale linearly with CPU frequency [18]. We plan to model frequency-independent parts in our future work.



Based on (6), a DRE system with  $m$  tasks and  $n$  processors is described by the following MIMO dynamic model.

$$\mathbf{u}(k) = \mathbf{u}(k-1) + \mathbf{G}\Delta\mathbf{b}(k-1) \quad (8)$$

where  $\mathbf{G}$  is a diagonal matrix where  $g_{ii} = g_i$  ( $1 \leq i \leq n$ ) and  $g_{ij} = 0$  ( $i \neq j$ ).  $\Delta\mathbf{b}(k)$  is a vector including the estimated utilization change (7) of each processor.

## 4.2 Controller Design

In this paper, we adopt the EUCON algorithm presented in our previous work [4] for rate adaptation. EUCON features a Model Predictive Controller (MPC) that optimizes a *cost function* defined over  $P$  control periods in the future, called the *prediction horizon*. The control objective is to select control inputs in the following  $M$  control periods, called *control horizon*, that minimize the following cost function while satisfying the constraints.

$$V(k) = \sum_{i=1}^P \|\mathbf{u}(k+i|k) - \mathbf{ref}(k+i|k)\|^2 + \sum_{i=0}^{M-1} \|\Delta\mathbf{r}(k+i|k) - \Delta\mathbf{r}(k+i-1|k)\|^2 \quad (9)$$

where  $P$  is the prediction horizon, and  $M$  is the control horizon. The first term in the cost function represents the *tracking error*, i.e., the difference between the utilization vector  $\mathbf{u}(k+i|k)$  and a reference trajectory  $\mathbf{ref}(k+i|k)$  defined in [4]. By minimizing the tracking error, the closed-loop system will converge to the utilization set points if the system is stable. The second term in the cost function represents the control penalty. This control problem is subject to the rate constraints (2). The detailed design and analysis of EUCON are available in [4].

Although the rate adaptation loop has been proved to be stable in [4], in order for the coordinated control architecture to be stable, the stability and settling time of the rate adaptation loop need to be reexamined by considering the impact from the frequency scaling loop. The detailed coordination analysis is presented in Section 5.4.

## 5 CPU Frequency Scaling Loop

In this section, we first model, design, and analyze the CPU frequency scaling loop. We then analyze the coordination between the two control loops.

### 5.1 System Model

Based on our control architecture, the frequency scaling loop can be designed separately from rate adaptation. As a result, model (6) can be simplified by having  $r_i(k)$  in (5) as a constant  $r_i$ . This decouples different processors because, as discussed in Section 3.1, processors are coupled to each other due to the fact that the rate change of a task may affect the utilizations of all the processors where its subtasks are located. The utilization of each processor can now be modeled individually because the CPU frequency change  $\Delta d_i(k) = 1/f_i(k) - 1/f_i(k-1)$  only affects the execution

times of all the subtasks on  $P_i$ . Specifically, the model of processor  $P_i$  is:

$$u_i(k) = u_i(k-1) + g_i \Delta d_i(k) \sum_{T_{jl} \in S_i} c_{jl} r_j \quad (10)$$

The model cannot be directly used to design controller because the system gain  $g_i$  is used to model the uncertainties in task execution times and thus unknown at design time. Therefore, we design the controller based on an approximate system model, which is model (10) with  $g_i = 1$ . In a real system where the task execution times are different than their estimations, the *actual* value of  $g_i$  may become different than 1. As a result, the closed-loop system may behave differently. However, in Section 5.3, we show that a system controlled by the controller designed with  $g_i = 1$  can remain stable as long as the variation of  $g_i$  is within a certain range. This range is established using stability analysis of the closed-loop system by considering the model variations.

### 5.2 Controller Design

Following standard control theory [19], we design a Proportional (P) controller to achieve the desired control performance such as stability and zero steady state error. We choose to use a P controller instead of a more sophisticated controller such as a PID (Proportional-Integral-Derivative) controller because the actuator  $1/f_i(k) = \Delta d_i(k) + 1/f_i(k-1)$  already includes an integrator such that zero steady state error can be achieved without resorting to an I (Integral) part. The D (Derivative) part is not used because it may amplify the noise in utilization in unpredictable environments. The Z-domain form of our P controller is:

$$C(z) = \frac{1}{\sum_{T_{jl} \in S_i} c_{jl} r_j} \quad (11)$$

The transfer function of the closed-loop system controlled by controller (11) is:

$$G(z) = z^{-1} \quad (12)$$

It is easy to prove that the controlled system is stable and has zero steady state errors when  $g_i = 1$ . The detailed proofs can be found in a standard control textbook [19] and are skipped due to space limitations. The desired CPU frequency in the  $k^{\text{th}}$  control period is:

$$f_i(k) = \frac{f_i(k-1) \sum_{T_{jl} \in S_i} c_{jl} r_j}{(U_s - u(k)) f_i(k-1) + \sum_{T_{jl} \in S_i} c_{jl} r_j} \quad (13)$$

### 5.3 Control Analysis for Model Variation

In this subsection, we analyze the system stability when the designed P controller is used on a system with  $g_i \neq 1$ . A fundamental benefit of the control-theoretic approach is that it gives us theoretical confidence for system stability, even when the task execution times are significantly different from their estimations.

The closed-loop transfer function for the real system is

$$G(z) = \frac{g_i}{z - (1 - g_i)} \quad (14)$$

The closed-loop system pole in (14) is  $1 - g_i$ . In order for the system to be stable, the pole must be within the unit circle. Hence, the system will remain stable as long as  $0 < g_i < 2$ . The result means that the actual utilization change *cannot* be twice the estimated utilization change.

We now analyze the steady state error of the controlled system when  $g_i \neq 1$ .

$$\lim_{z \rightarrow 1} (z - 1)U(z) = \lim_{z \rightarrow 1} \left( \frac{g_i z}{z - (1 - g_i)} U_s \right) = U_s \quad (15)$$

Equation (15) means that we are guaranteed to achieve the desired CPU utilization as long as the system is stable.

## 5.4 Coordination Analysis

We now analyze the coordination needed for the two control loops to work together with global stability. The analysis here, as well as the control architecture design in Section 3 and our empirical results, demonstrates the importance of coordinating different control loops, which is a major contribution of our paper.

First, we need to ensure that the stability of the rate adaptation loop will not be affected when the frequency scaling loop changes the CPU frequency and so  $f_i(k) \neq 1$ . Given a specific task set, the stability condition of the rate adaptation loop as a range of  $g_i$  (i.e., the ratio between the actual utilization change and the estimated change) can be established by following the steps presented in [4]. For example, the stability condition of the task set used in our experiments is that the actual change cannot be 10 times the estimated change. Accordingly, we must guarantee that the relative CPU frequency of each processor is not smaller than 0.1 because the rate adaptation controller is designed with the assumption of  $f_i(k) = 1$ . This constraint must be enforced in the frequency scaling loop. One of the reasons for us to choose the rate adaptation loop as the secondary loop in our control architecture is that it has a larger stability range and thus is less sensitive to the impact of the primary loop.

Second, we must guarantee that the frequency scaling loop is also stable, i.e.,  $0 < g_i < 2$ . Since the frequency scaling loop is the primary loop of our two-layer control architecture, the difference between the actual and estimated utilization changes is mainly caused by the differences between the actual and estimated execution times. Therefore, it is preferable to use pessimistic estimation on execution times such that the controlled system can be guaranteed to be stable and the system oscillation can also be reduced. Please note that using pessimistic estimated execution times does not result in underutilization of the CPU as in systems that rely on traditional open-loop scheduling. This is because our control architecture dynamically adjusts CPU frequencies and tasks rates based on *measured* utilization

rather than the estimated execution times. The downside of using more pessimistic estimation on execution times is that it leads to a smaller system gain, which may cause slower convergence to the set points. However, since it is more important to guarantee system stability in a DRE system, it is still preferable to overestimate task execution times.

Third, we need to analyze the settling time of the rate adaptation loop in order to determine the control period of the frequency scaling loop. Since settling time has not been analyzed in [4], we now outline the general process of analyzing the settling time of the rate adaptation loop when the actual utilization change is different from the estimated change, i.e.,  $g_i \neq 1$ . First, given a specific task set, we derive the control inputs  $\Delta \mathbf{r}(k)$  that minimize the cost function (9) based on the system model (8) with  $g_i = 1$ . The control inputs represent the control decision based on the estimated system model. Second, we derive the closed-loop system model by substituting the control inputs derived in the first step into the system model (8) where  $g_i \neq 1$ . The analysis needs to consider a composite system consisting of the dynamics of the original system and the controller. Finally, we calculate the dominant pole (i.e., the pole with the largest magnitude) of the closed-loop system. According to control theory, the dominant pole determines the system's transient response such as settling time.

Based on our analysis, the task set used in our experiments has a settling time of 5 control periods under rate adaptation. The detailed derivation is not included due to space limitations. The control period of the rate adaptation loop is selected to be 2 seconds to include multiple instances of each task, resulting in a settling time of  $10 = 5 \times 2$  seconds. Therefore, the control period of the frequency scaling loop is set to 20 seconds, which is much longer than the settling time of the rate adaptation loop.

## 6 System Implementation

Our testbed includes 4 Linux servers, called RTES1 to RTES4, to run the end-to-end real-time tasks and a desktop machine to run the MPC controller. The 4 servers are equipped with 2.4GHz AMD Athlon 64 3800+ processors with 1GB RAM and 512KB L2 Cache. The controller machine is a Dell OptiPlex GX520 with 3.00GHz Intel Pentium D Processor and 1GB RAM. All the machines are connected by a 100Mbps internal Ethernet switch. The 4 servers run openSUSE Linux 11 with kernel 2.6.25 while the controller machine runs Windows XP.

We implement our control architecture in FC-ORB, an open-source real-time Object Request Broker (ORB) middleware system [20]. FC-ORB supports end-to-end real-time tasks based on the end-to-end scheduling framework [12]. FC-ORB implements the release guard protocol to enforce the precedence constraints among subtasks.

Our experiments run a medium-sized workload that comprises 12 end-to-end tasks (with a total of 25 subtasks). The subtasks on each processor are scheduled by the RMS algorithm [12]. Each task's end-to-end deadline

is  $d_i = n_i/r_i(k)$ , where  $n_i$  is the number of subtasks in task  $T_i$  and  $r_i(k)$  is the current rate of  $T_i$ . Each end-to-end deadline is evenly divided into subdeadlines for its subtasks. The resultant subdeadline of each subtask  $T_{ij}$  equals its period,  $1/r_i(k)$ . The utilization set point of every processor is set to its RMS schedulable utilization bound [12], i.e.,  $B_i = n_i(2^{1/n_i} - 1)$ , where  $n_i$  is the number of subtasks on  $P_i$ . All (sub)tasks meet their (sub)deadlines if the desired utilization on every processor is enforced.

We now introduce the implementation details of each component in our two-layer control architecture.

**Utilization Monitor:** The utilization monitor uses the `/proc/stat` file in Linux to estimate the CPU utilization in each control period. The `/proc/stat` file records the number of jiffies (usually 10ms in Linux) when the CPU is in user mode, user mode with low priority (`nice`), system mode, and when used by the idle task, since the system starts. At the end of each control period, the utilization monitor reads the counters, and estimates the CPU utilization as 1 minus the number of jiffies used by the idle task in the last control period and then divided by the total number of jiffies in the same period.

**MPC Controller:** The controller is implemented as a single-thread process running separately on the controller machine. Each time its periodic timer fires, the controller sends utilization requests to all the 4 application servers. The incoming replies are handled asynchronously so that the controller can avoid being blocked by an overloaded application server. After the controller collects the replies from all the servers, it executes the control algorithm introduced in Section 4.2 to calculate the new task rates. The controller then sends the tasks' new rates to the rate modulators on the servers for enforcement. If a server does not reply in an entire control period, its utilization is treated as 100%, as the controller assumes this server is overloaded with its (sub)tasks and so cannot respond. The control period of the rate adaptation loop is 2 seconds.

**Rate Modulator:** A Rate Modulator is located on each processor. It receives the new rates from the controller and then resets the timer interval of the first subtask of each task whose invocation rate needs to be changed.

**Proportional Controller:** The controller is implemented as a process running on each of the 4 servers. With a control period of 20 seconds, the controller periodically reads the CPU utilization of the server, executes the control algorithm presented in Section 5.2 to compute the desired CPU frequency, and sends the new frequency to the frequency modulator on the server.

**Frequency Modulator:** We use AMD's Cool'n'Quiet technology to enforce the new CPU frequency. AMD Athlon 64 3800+ microprocessor has 5 discrete CPU frequency levels. To change CPU frequency, one needs to install the `cpufreq` package and then use root privilege to write the new frequency level into the system file `/sys/devices/system/cpu/cpu0/cpufreq/scaling_setspeed`. A routine periodically checks this file and resets the CPU frequency

accordingly. The average overhead (i.e., transition latency) to change frequency in AMD Athlon processors is about  $100\mu s$  according to the AMD white paper report.

Since the new CPU frequency level periodically received from the proportional controller could be any value that is not exactly one of the five supported frequency levels. Therefore, the modulator code must locally resolve the output value of the controller to a series of supported frequency levels to approximate the desired value. For example, to approximate 2.89GHz during a control period, the modulator would output the sequence 2.67, 3, 3, 2.67, 3, 3, etc on a smaller timescale. To do this, we implement a first-order delta-sigma modulator, which is commonly used in analog-to-digital signal conversion. The detailed algorithm of the first-order delta-sigma modulator can be found in [15].

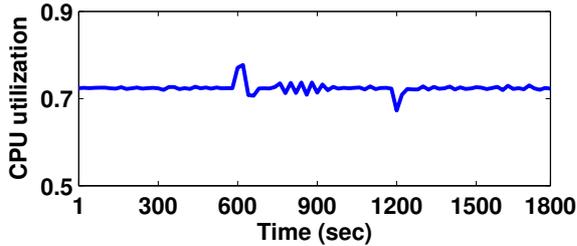
**Power Monitor:** The power consumption of each server is measured with a WattsUp Pro power meter by plugging the server into the power meter, which is connected to a standard 120V AC wall outlet. The WattsUp power meter has an accuracy of  $\pm 1.5\%$  of the measured value. To access power data, the data port of each power meter is connected to a serial port of the data collection machine. The power meter samples the power data every second and then sends the reading to the data collection program through a system file `/dev/ttyUSB0`.

## 7 Empirical Results

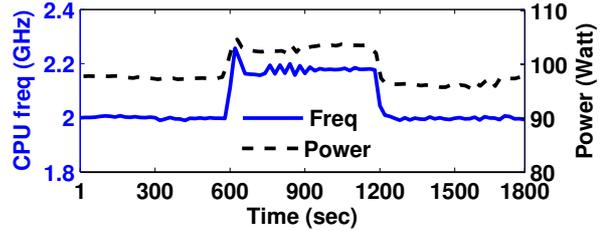
In this section, we first test the frequency scaling loop alone. We then show that the frequency scaling loop can effectively control utilizations when it is infeasible for a rate adaptation controller to do so. Finally, we demonstrate that the coordinated control solution can maximize the system's adaptation capability for power-efficient utilization control.

### 7.1 Frequency Scaling Loop

In this experiment, we disable the rate adaptation loop to evaluate the performance of the frequency scaling loop on server RTES1. As a common practice in real-time systems that rely on open-loop scheduling algorithms, the workload of RTES1 is configured with carefully tuned initial task rates such that the server has an initial CPU utilization of 0.72, which is its RMS bound. As shown in Figure 2(a), at time 600s, the execution times of all the tasks on RTES1 are suddenly increased by 8% to test the system's capability of handling workload fluctuations. The increase makes the CPU utilization of RTES1 jump to 0.78, which is higher than the RMS bound and so may cause undesired deadline misses. Figure 2(b) shows that the frequency scaling loop responds to the utilization increase by dynamically increasing the CPU frequency of the server processor from 2.0GHz to 2.18GHz. As a result, the utilization returns back to the set point quickly. In contrast, An open-loop system without dynamic feedback would have its utilization stay above the RMS bound. At time 1200s, the task execution times are suddenly reduced back to their original values, resulting in



(a) CPU utilization



(b) CPU frequency and power consumption

**Figure 2. CPU utilization control by frequency scaling under a workload increase from 600s to 1200s.**

a utilization lower than the set point. The frequency scaling loop then responds by reducing the CPU frequency back to 2.0GHz for power savings.

To test the robustness of the controller, we conduct a set of experiments with different utilization set points. Figure 3(a) plots the means and the standard deviations of RTES1’s CPU utilization after the controller enters the steady state. We can see that the frequency scaling loop can successfully achieve the desired utilization set points. Figure 3(b) demonstrates that more power saving has been achieved when we allow the system to have a utilization set point closer to its RMS schedulable bound, i.e., 0.72. The maximum standard deviation for power is smaller than 1.5 W.

## 7.2 Frequency Scaling vs. EUCON

In this experiment, we show that frequency scaling can be used to control CPU utilizations when rate adaptation fails to do so in some cases. We compare the frequency scaling loop with a baseline, a state-of-the-art control algorithm called EUCON [4], which relies only on the rate adaptation loop briefly introduced in Section 4. Figure 4(a) shows that EUCON fails to achieve the desired set points (0.74 for RTES2 and 0.72 for the other three servers) because the task rates saturate at the upper boundaries of their allowed ranges. As a result, the system is underutilized with unnecessarily high power consumption, as shown in Figure 5(a). We then test the frequency scaling loop using the same workload with the rate adaptation loop disabled. In the experiment, to highlight the performance of the frequency scaling loop, we first let the system run in an open-loop manner (with no controller activated). Therefore, the system initially cannot achieve the desired CPU utilizations. At time 400s, we activate the frequency scaling loop. Figure 4(b) shows that the CPU utilizations quickly converge to their desired set points. As a result, all the servers achieve power savings (as shown in Figure 5(b)) while still guaranteeing the end-to-end task schedulability.

## 7.3 Coordinated Utilization Control

Since both task rates and CPU frequencies can only be adapted within allowed ranges, our coordinated control solution is designed to combine them based on control theory for maximized adaptation capability. In this experiment, we run the same workload with all the tasks starting with lower initial rates than those used in Section 7.2. As a result, Figure 6(a) shows that the utilizations controlled by the rate

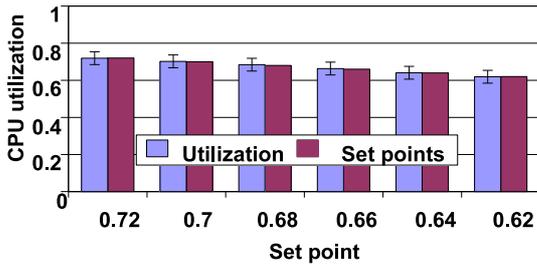
adaptation loop start from values lower than those in Figure 4(a). Similar to Figure 4(a), the rate adaptation loop fails to achieve the desired utilization set points (dashed lines in the figure) because tasks are already running at their highest possible rates allowed by their ranges. In this case, the CPU frequencies of the processors could be lowered for power savings. We then examine the frequency scaling loop alone by running the same experiment in Section 7.2 with lower initial task rates. Figure 4(b) shows that the frequency scaling loop fails to achieve the desired utilizations this time because the tasks are running at lower rates. As a result, even when the processors are already running at their lowest CPU frequencies, utilizations still cannot converge to the desired set points. In this case, we could allow tasks to run at higher rates to contribute a higher value to the system.

We now evaluate our coordinated control solution. To highlight the performance of our solution, we first run the rate adaptation loop, which achieves the highest rates for all the tasks, resulting in a high system value. At time 420s, we activate the frequency scaling loop. Figure 7(a) shows that the coordinated control solution successfully achieves the desired utilization set points. In the meantime, Figure 7(b) demonstrates that servers RTES2, RTES3, RTES4 also receive considerable power savings. Therefore, the coordinated control solution can effectively control CPU utilizations to the desired set points while achieving increased task rates and reduced power consumption.

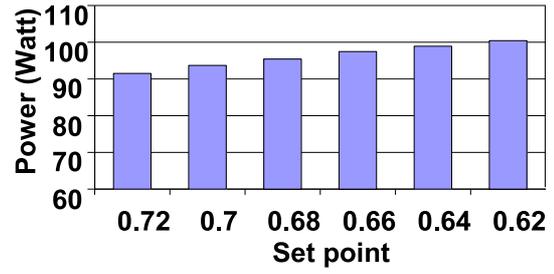
## 8 Related Work

A survey of feedback performance control in computing systems is presented in [21]. Many projects that applied control theory to real-time scheduling and applications are closely related to this paper. Steere et al. and Goel et al. developed feedback-based schedulers [22][23] that guarantee desired progress rates for real-time applications. Abeni et al. presented control analysis of a reservation-based feedback scheduler [24]. Lu et al. developed a middleware service that adopts feedback control scheduling algorithms to control CPU utilization and deadline miss ratio [6]. Feedback control has also been applied to power control [25][15] and digital control applications [26].

Various CPU utilization control algorithms (e.g., [6][27][28][20]) have been recently proposed to guarantee real-time deadlines. For example, Lu et al. designed constrained MIMO utilization control algorithm for multi-

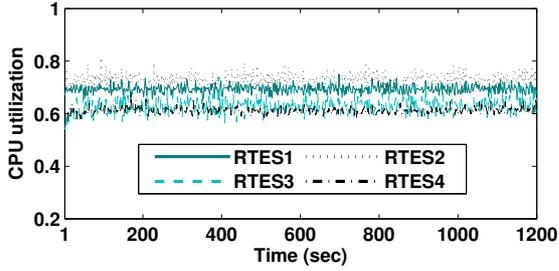


(a) CPU utilization

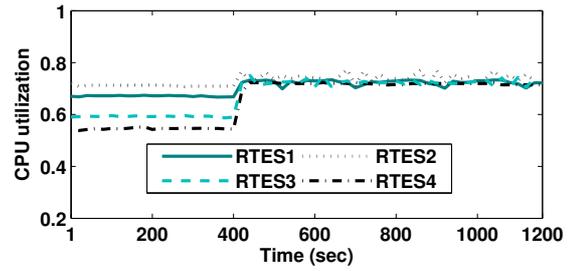


(b) Power consumption

**Figure 3. CPU utilization control by frequency scaling under different utilization set points.**

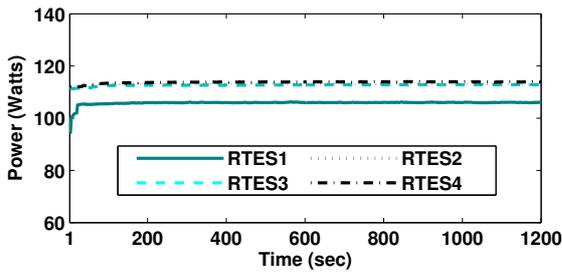


(a) EUCON

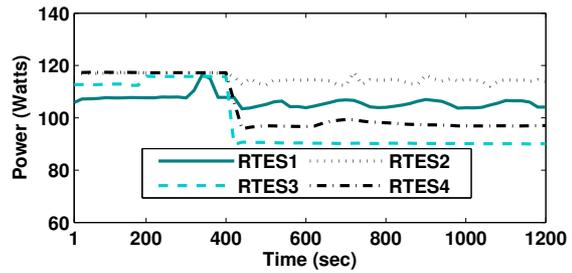


(b) Frequency scaling activated from 400s

**Figure 4. Comparison of control accuracy between EUCON and the frequency scaling loop.**

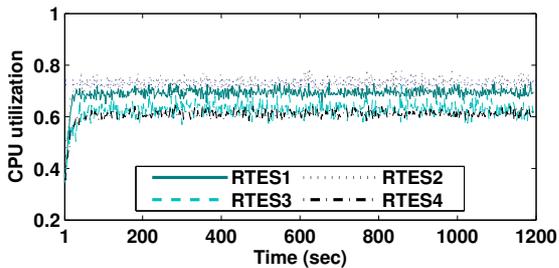


(a) EUCON

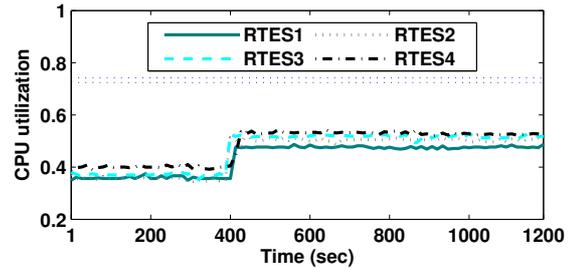


(b) Frequency scaling activated from 400s

**Figure 5. Comparison of power consumption between EUCON and the frequency scaling loop.**

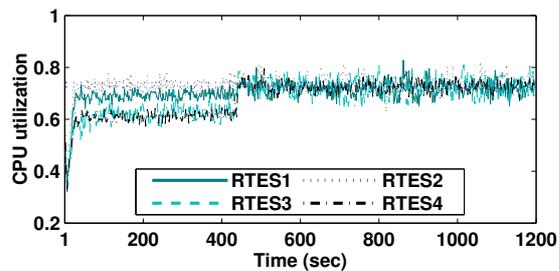


(a) Rate adaptation

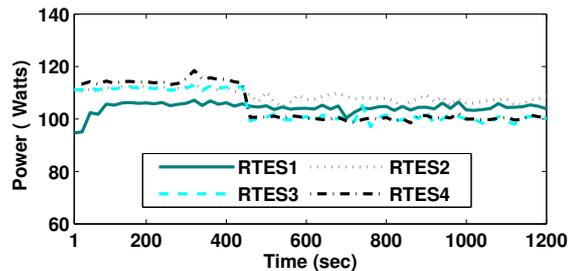


(b) Frequency scaling activated from 400s

**Figure 6. Infeasible utilization control by rate adaptation or frequency scaling individually.**



(a) CPU utilization



(b) Power consumption

**Figure 7. CPU utilization control by the coordinated control solution from 420s.**

ple processors that are coupled due to end-to-end tasks [4]. Wang et al. proposed decentralized utilization control algorithm for large-scale distributed real-time systems [5]. Yao et al. developed an adaptive utilization control algorithm [29]. However, all those algorithms assume that task rates can only be continuously tuned. Hybrid control theory [9] and optimization strategies [8] are adopted to handle discrete task rates based on the assumption that task WCETs are known *a priori* and accurate. In contrast to all the existing work that relies exclusively on rate adaptation, we present a two-layer control architecture that uses both rate adaptation and DVFS for power-efficient utilization control.

Many energy-efficient real-time scheduling algorithms have been proposed (e.g., [18][30][31][17][16][25]). Most existing work relies on detailed knowledge (e.g., WCETs) of workloads to minimize the energy consumption or temperature, or maximize the system reward in an *open-loop* manner. While they can effectively guarantee task schedulability in closed environments without a feedback loop for adaptations, they may not be directly applied to DRE systems whose workloads may vary significantly at runtime. In contrast, we use DVFS as a knob to dynamically react to unpredictable workload variations instead of minimizing the energy consumption of the entire DRE system.

## 9 Conclusions

In this paper, we have formulated a new CPU utilization control problem based on both frequency scaling and rate adaptation. Since a centralized controller for simultaneous frequency scaling and rate adaptation would have a non-linear system model, we designed a two-layer coordinated CPU utilization control architecture. The primary control loop uses frequency scaling to locally control the CPU utilization of each processor, while the secondary control loop adopts rate adaptation to control the utilizations of all the processors in the system at the cluster level on a smaller timescale. Both the two control loops are designed and coordinated based on well-established control theory for theoretically guaranteed control accuracy and global system stability. Empirical results on a physical testbed demonstrate that our control solution outperforms EUCON, a state-of-the-art utilization control algorithm, by having increased adaptation capability and less power consumption.

## References

- [1] J. Sun and J. Liu, "Synchronization protocols in distributed real-time systems," in *ICDCS*, 1996.
- [2] R. Rajkumar, L. Sha, and J. P. Lehoczky, "Real-time synchronization protocols for multiprocessors," in *RTAS*, 1988.
- [3] D. Henriksson and T. Olsson, "Maximizing the use of computational resources in multi-camera feedback control," in *RTAS*, 2004.
- [4] C. Lu, X. Wang, and X. Koutsoukos, "Feedback utilization control in distributed real-time systems with end-to-end tasks," *IEEE Transactions on Parallel and Distributed Systems*, vol. 16, no. 6, 2005.
- [5] X. Wang, D. Jia, C. Lu, and X. Koutsoukos, "DEUCON: Decentralized end-to-end utilization control for distributed real-time systems," *IEEE Trans. on Parallel and Distributed Sys.*, vol. 18, no. 7, 2007.
- [6] C. Lu, X. Wang, and C. Gill, "Feedback control real-time scheduling in ORB middleware," in *RTAS*, 2003.
- [7] X. Wang, Y. Chen, C. Lu, and X. Koutsoukos, "On controllability and feasibility of utilization control in distributed real-time systems," in *ECRTS*, 2007.
- [8] Y. Chen, C. Lu, and X. Koutsoukos, "Optimal discrete rate adaptation for distributed real-time systems," in *RTSS*, 2007.
- [9] X. Koutsoukos, R. Tekumalla, B. Natarajan, and C. Lu, "Hybrid supervisory utilization control of real-time systems," in *RTAS*, 2005.
- [10] S. Goddard and X. Liu, "A variable rate execution model," in *ECRTS*, 2004.
- [11] S. A. Brandt and G. J. Nutt, "Flexible soft real-time processing in middleware," in *RTSS*, 2004.
- [12] J. W. S. Liu, *Real-Time Systems*. Prentice Hall, 2000.
- [13] P. Marti, G. Fohler, P. Fuertes, and K. Ramamritham, "Improving quality-of-control using flexible timing constraints: metric and scheduling," in *RTSS*, 2002.
- [14] G. C. Buttazzo, G. Lipari, M. Caccamo, and L. Abeni, "Elastic scheduling for flexible workload management," *IEEE Transactions on Computers*, vol. 51, no. 3, 2002.
- [15] C. Lefurgy, X. Wang, and M. Ware, "Power capping: a prelude to power shifting," *Cluster Computing*, vol. 11, no. 2, 2008.
- [16] H. Aydin, P. Mejía-Alvarez, D. Mossé, and R. Melhem, "Dynamic and aggressive scheduling techniques for power-aware real-time systems," in *RTSS*, 2001.
- [17] S. Saewong and R. R. Rajkumar, "Practical voltage-scaling for fixed-priority RT-systems," in *RTAS*, 2003.
- [18] H. Aydin, V. Devadas, and D. Zhu, "System-level energy management for periodic real-time tasks," in *RTSS*, 2006.
- [19] G. F. Franklin, J. D. Powell, and M. Workman, *Digital Control of Dynamic Systems, 3rd edition*. Addison-Wesley, 1997.
- [20] X. Wang, Y. Chen, C. Lu, and X. Koutsoukos, "FC-ORB: A robust distributed real-time embedded middleware with end-to-end utilization control," *Journal of Systems and Software*, vol. 80, no. 7, 2007.
- [21] T. F. Abdelzaher, J. Stankovic, C. Lu, R. Zhang, and Y. Lu, "Feedback performance control in software services," *IEEE Control Systems*, vol. 23, no. 3, 2003.
- [22] D. C. Steere, A. Goel, J. Gruenberg, D. McNamee, C. Pu, and J. Walpole, "A feedback-driven proportion allocator for real-rate scheduling," in *OSDI*, 1999.
- [23] A. Goel, J. Walpole, and M. Shor, "Real-rate scheduling," in *RTAS*, 2004.
- [24] L. Abeni, L. Palopoli, G. Lipari, and J. Walpole, "Analysis of a reservation-based feedback scheduler," in *RTSS*, 2002.
- [25] Y. Zhu and F. Mueller, "Feedback EDF scheduling exploiting dynamic voltage scaling," in *RTAS*, 2004.
- [26] A. Cervin, J. Eker, B. Bernhardsson, and K.-E. Arzen, "Feedback-feedforward scheduling of control tasks," *Real-Time Systems*, vol. 23, no. 1, 2002.
- [27] J. A. Stankovic, T. He, T. Abdelzaher, M. Marley, G. Tao, S. Son, and C. Lu, "Feedback control scheduling in distributed real-time systems," in *RTSS*, 2001.
- [28] S. Lin and G. Manimaran, "Double-loop feedback-based scheduling approach for distributed real-time systems," in *HiPC*, 2003.
- [29] J. Yao, X. Liu, M. Yuan, and Z. Gu, "Online adaptive utilization control for real-time embedded multiprocessor systems," in *CODES+ISSS*, 2008.
- [30] R. Xu, R. Melhem, and D. Moss, "Energy-aware scheduling for streaming applications on chip multiprocessors," in *RTSS*, 2007.
- [31] J.-J. Chen, C.-M. Hung, and T.-W. Kuo, "On the minimization of the instantaneous temperature for periodic real-time tasks," in *RTAS*, 2007.