

# Optimal Discrete Rate Adaptation for Distributed Real-Time Systems<sup>\*</sup>

Yingming Chen   Chenyang Lu   Xenofon Koutsoukos  
Washington University in St. Louis   Vanderbilt University

## Abstract

Many distributed real-time systems face the challenge of dynamically maximizing system utility and meeting stringent resource constraints in response to fluctuations in system workload. Thus, online adaptation must be adopted in face of workload changes in such systems. We present the MultiParametric Rate Adaptation (MPRA) algorithm for discrete rate adaptation in distributed real-time systems with end-to-end tasks. The key novelty and advantage of MPRA is that it can efficiently produce optimal solutions in response to workload variations such as dynamic task arrivals. Through offline preprocessing MPRA transforms an NP-hard utility optimization problem to the evaluation of a piecewise linear function of the CPU utilization. At run time MPRA produces optimal solutions by evaluating the function based on the CPU utilization. Analysis and simulation results show that MPRA maximizes system utility in the presence of varying workloads, while reducing the online computation complexity to polynomial time.

## 1 Introduction

An increasing number of distributed real-time systems operate in dynamic environments where system workload may change at run time [1]. A key challenge faced by such systems is to dynamically maximize system utility subject to resource constraints and fluctuating workload. For instance, the Supervisory Control and Data Acquisition (SCADA) system of a power grid may experience dramatic load increase during cascading power failures and cyber attacks. Similarly, the arrival rate of service requests in an online trading server can fluctuate dramatically. However, such systems must meet stringent resource constraints despite their fluctuating workload. In particular, such systems need to enforce desired CPU utilization bounds on multiple processors in order to provide overload protection and meet end-to-end deadlines. Therefore, online adaptation must be adopted to handle workload changes in such systems.

Online adaptation introduces several important challenges. First, online adaptation should *maximize system utility* subject to multiple resource constraints. For example, many distributed real-time systems must enforce certain CPU utilization bounds on multiple processors in order to prevent system crash due to CPU saturation and meet end-to-end deadlines. Second, many common adaptation strategies only support *discrete* options. For example, an admission controller must make binary decisions (admission/rejection) on a task. While task rate adaptation can allow a system to adapt at a finer granularity [8][10][23][15][30], many real-time applications (e.g., avionics [2] and Multiple Bit-Rate Video) can only run at a discrete set of predefined rates. Unfortunately, utility optimization problems with discrete options are NP-hard [18]. Furthermore, despite the difficulty of such problems, a real-time system must adapt to dynamic workload changes quickly, which requires optimization algorithms to be highly efficient at run time.

Existing approaches to utility optimization in real-time systems can be divided into two categories: optimal solutions and efficient heuristics. Approaches based on integer programming or dynamic programming have been proposed to optimize utility [18][17]. While these approaches produce optimal solutions, they are computationally expensive and cannot be used *online*. On the other hand, a number of efficient heuristics have been proposed for online adaptation [28][18][2][20]. However, these algorithms can only produce sub-optimal solutions in terms of system utility.

To overcome the limitations of existing approaches, we present the *MultiParametric Rate Adaptation (MPRA)* algorithm for online adaptation in real-time systems. MPRA employs task rate adaptation as the online adaptation mechanism, which is supported by a broad range of real-time applications, such as digital control [10], video streaming, and avionics [2]. Specifically, MPRA is designed to handle end-to-end tasks that may only execute at a discrete set of rates on multiple processors. This task model introduces significant challenges to optimal online adaptation algorithms.

The key novelty and advantage of our approach is that it can *efficiently* produce *optimal* solutions online in face of workload changes caused by dynamic task arrivals and departures. The MPRA algorithm is based on multipara-

<sup>\*</sup>This research was supported by NSF CAREER award (grant CNS-0448554).

metric mixed-integer linear programming (mp-MILP) [4]. Through offline preprocessing MPRA transforms an NP-hard utility optimization problem to the evaluation of a piecewise linear function of the CPU utilization. At run time MPRA produces optimal solutions by evaluating the function based on the workload variation. Specifically, the primary contributions of this paper are three-fold:

- We present MPRA, a novel algorithm for discrete rate adaptation in distributed real-time systems with end-to-end tasks;
- We provide analysis that proves that our algorithm produces optimal system utility in face of workload changes with the online rate adaptation running in *polynomial* time;
- We present simulation results that demonstrate that MPRA maximizes system utility in the presence of dynamic task arrivals, with the online execution time comparable to efficient suboptimal heuristics and two orders of magnitude lower than a representative optimal solver.

The rest of the paper is organized as follows. Section 2 discusses related work. Section 3 formalizes the optimization problem addressed in this paper. Section 4 presents the design and complexity analysis of our algorithm. Section 5 provides simulation results. Finally, Section 6 concludes this paper.

## 2 Related Work

Several projects investigated the problem of maximizing system utility in real-time systems. Rajkumar et al. proposed the QoS-based Resource Allocation Model (Q-RAM) [27] for utility optimization in distributed real-time systems. Lee et al. presented several optimal algorithms for the Q-RAM model based on integer programming and dynamic programming [18][17]. These approaches are computationally expensive and unsuitable for online adaptation in real-time systems. To improve the efficiency of the solutions, the authors also proposed several efficient heuristic algorithms that can only produce sub-optimal solutions [28][18][17][19][12]. Specifically, they presented heuristic algorithms with bounded approximation ratio for the single-resource case [17][19]. However, the heuristic algorithms for multi-resource problems do not have analytical bounds on the approximation ratio [19]. Note that the multi-resource case is common in distributed real-time systems in which each processor is a separate resource.

Various system-wide schemes have been studied to improve system utility. The authors in [32][6][7] have developed middleware solutions that support mediating applica-

tion resource usage using application QoS levels for single processor systems. Abdelzaher et al. developed a QoS-negotiation model and incorporated it into an example real-time middleware service, called RTPOOL, in [2]. All of the projects developed middleware systems that aim to improve system utility by dynamically adjusting the QoS levels of applications. However, they employ heuristic algorithms that cannot produce optimal solutions.

Recently, Lee et al. introduced a method called service class configuration to address the online adaptation problem with dynamic arrival and departure of tasks in distributed real-time systems [20]. This method avoids running optimization procedures at run time by designing a set of service classes offline, which will be used adaptively depending on the system state. While service classes can effectively improve the efficiency of online adaptation, it cannot produce optimal solutions. In contrast, MPRA can produce optimal solutions with efficient online execution.

Several task rate adaptation algorithms have been proposed for single-processor [10][8][30] and distributed real-time systems [23][35]. All the above solutions assume that task rates can be adjusted in a *continuous* range. As discussed in Section 1, this assumption does not hold in many applications that only support *discrete* configurations. HySUCON [15] is a heuristic algorithm for real-time systems that supports discrete task rates. However, it is designed for single processor systems and cannot produce optimal solutions. There are two important differences between our work and earlier work on rate adaptation. First, our work deals with real-time systems with *discrete* task rates, while none of the aforementioned rate adaptation algorithms (with the exception of [15]) is designed to handle discrete rates. Moreover, none of them can maximize system utility.

## 3 Problem Formulation

We now formulate the discrete rate adaptation problem in distributed real-time systems.

### 3.1 End-to-End Task Model

We classify tasks in distributed real-time systems into two categories: *adaptable* tasks and *unadaptable* tasks. Tasks that support multiple rate choices are called *adaptable* tasks; tasks with fixed rates fall into *unadaptable* task category. Mission critical tasks that must execute at fixed rates are typical *unadaptable* tasks.

We assume the system is comprised of  $m$  *adaptable* 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 graph of subtasks  $\{T_{ij} | 1 \leq j \leq m_i\}$  that may be located on different processors. We denote the set of subtasks of *adaptable* task

$T_i$  that are allocated on  $P_j$  as  $S_{ji}$ . Due to the dependencies among subtasks each subtask  $T_{ij}$  of a periodic task  $T_i$  shares the same rate as  $T_i$ .<sup>1</sup> Each task  $T_i$  is subject to an end-to-end relative deadline related to its period  $\tau_i$ . Each subtask  $T_{ij}$  has an execution time  $c_{ij}$ .

We assume each *adaptable* task only supports a set of discrete task rates for online adaptation. A task running at a higher rate contributes a higher utility to the system at the cost of higher utilization. We denote the set of discrete rate choices of task  $T_i$  as  $R_i = \{r_i^{(0)}, \dots, r_i^{(k_i)}\}$  in increasing order. The set of utility options for task  $T_i$  is denoted by  $Q_i = \{q_i^{(0)}, \dots, q_i^{(k_i)}\}$  where  $q_i^{(j)}$  is the utility value contributed by  $T_i$  when it is configured with  $r_i^{(j)}$ . Note that we do not make any assumption regarding the relation between the task utility and the task rate. For example, a task's utility values do not need to be a linear or polynomial function of the task rate. MPRA can handle arbitrary utility values assigned to discrete task rates. Task utility values for different rates can be represented by a lookup table, which is specified by application designers based on domain knowledge. Admission control is a special case of discrete rate adaptation, in which each task only have two rate choices: zero when the task is evicted and a fixed non-zero rate when task is admitted.

### 3.2 Discrete Rate Adaption Problem

Before formulating the discrete rate adaptation problem, we first introduce several notations:

- $R = [r_1, \dots, r_m]$  is the task rate vector where  $r_i$  is the current invocation rate of task  $T_i$ . Therefore we have  $r_i \in R_i, 1 \leq i \leq m$ .
- $Q_s$  is the system utility, i.e., the combined utility of all *adaptable* tasks defined as the weighted sum of the task utilities  $Q_s = \sum_{i=1}^m w_i q_i$  where  $q_i, q_i \in Q_i$ , is the current task utility of  $T_i$  and  $0 \leq w_i \leq 1, 1 \leq i \leq m$ , are weights describing the relative importance of the tasks.
- $D = [d_1, \dots, d_n]$  is the workload variation vector where  $d_i$  is the change to the utilization of the  $i^{th}$  processor caused by dynamic arrivals or departures of *unadaptable* tasks.  $D$  can be calculated based on the worst case execution times and rates of the *unadaptable* tasks that are assumed to be known. For example, denote the set of *unadaptable* tasks that arrive as  $S_a$  and the set of *unadaptable* tasks that just depart as  $S_b$ . Then  $d_i = \sum_{T_j \in S_a} \sum_{T_{jl} \in S_{ij}} c_{jl} r_j -$

$\sum_{T_j \in S_b} \sum_{T_{jl} \in S_{ij}} c_{jl} r_j$  where  $S_{ij}$  is the set of subtasks of  $T_j$  that run on processor  $P_i$ ,  $c_{jl}$  is the worst-case execution time of subtask  $T_{jl}$ , and  $r_j$  is the task rate of  $T_j$ .

- $U = [u_1, \dots, u_n]$  is the CPU utilization vector where  $u_i$  represents the utilization of the  $i^{th}$  processor.  $u_i$  is calculated by  $u_i = d_i + \sum_{1 \leq j \leq m} \sum_{T_{jl} \in S_{ij}} c_{jl} r_j$ .
- $B = [b_1, \dots, b_n]$  is the utilization bound vector where  $b_i$  is the utilization bound of the  $i^{th}$  processor specified by user.

The discrete rate adaptation problem can be formulated as a constrained optimization problem. The goal is to maximize the system utility via rate adaptation in response to workload changes when *unadaptable* tasks dynamically arrive or depart, i.e.

$$\max_R \sum_{i=1}^m w_i q_i \quad (1)$$

subject to

$$r_i \in R_i, 1 \leq i \leq m \quad (2)$$

$$U \leq B \quad (3)$$

The constraint (2) indicates that each task can only be configured with predefined rates. The utilization constraint (3) is used to enforce certain CPU utilization bounds on multiple processors in order to meet two real-time requirements:

*Meeting end-to-end deadlines.* Real-time tasks must meet their end-to-end deadlines in distributed real-time systems. In the end-to-end scheduling approach [31], the deadline of an end-to-end task is divided into subdeadlines of its subtasks, and the problem of meeting the end-to-end deadline is transformed to the problem of meeting the subdeadline of each subtask. A well-known approach for meeting the subdeadlines on a processor is by enforcing the schedulable utilization bound [22][21]. To meet end-to-end deadlines, the utilization set point of each processor is set to a value below its schedulable utilization bound. We can apply various subdeadline assignment algorithms [14][26] and schedulable utilization bounds for different task models [22][21] presented in the literature.

*Overload protection.* Many distributed systems must avoid saturation of processors, which may cause system crash or severe service degradation [3]. On COTS operating systems that support real-time priorities, high utilization by real-time threads may cause kernel starvation. The utilization constraint (3) allows a user to enforce desired utilization bounds for all the processors in a distributed system.

The discrete rate adaptation problem is NP-hard as it can be easily reduced to the 0-1 Knapsack Problem [24]. It is therefore impractical to apply standard optimization approaches to discrete rate adaptation in distributed real-time systems. There exist several approximation algorithms

<sup>1</sup>A non-greedy synchronization protocol [31] can be used to remove release jitter of subtasks.

for the 0-1 Knapsack Problem that run in polynomial time [13][29]. However, those algorithms can only handle problems for the single-resource case and can not be applied for multi-resource problems.

## 4 Design and Analysis of MPRA

In this section, we present the design and analysis of MPRA. The MPRA algorithm is based on Multiparametric programming, which is a general framework for solving mathematical programming problems with constraints that depend on varying parameters [11]. This technique is suitable for discrete rate adaptation problems, where the utilization constraints are related to online workload variations. It reduces the online utility optimization problem to a simple function evaluation problem. In the rest of this section, we first give a brief overview of the general framework of multiparametric programming. Next, we transform the discrete rate adaptation problem to an mp-MILP problem and design MPRA that instantiates the multiparametric programming approach for optimal and efficient rate adaptation in distributed real-time systems. Finally, we present the complexity analysis of our algorithm.

### 4.1 Multiparametric Programming

Multiparametric programming provides a systematic way to analyze the effect of parameter changes on the optimal solution of a mathematical programming problem. Rather than solving the optimization problem online, it includes an offline and an online algorithm. The offline algorithm partitions the space of varying parameters into regions. For each region, the objective and optimization variables are expressed as linear functions of the varying parameters. For a given value of the varying parameter, the online algorithm computes the optimal solution by evaluating the function for the region which includes the parameter value. The multiparametric approach has been extended for multiparametric mixed-integer linear programming problems (mp-MILP) [4]. The algorithm presented in [4] uses a Branch and Bound strategy to solve multiparametric 0-1 mixed-integer linear programming problems where the elements of the optimization vector can be either continuous or binary variables and parameters vary in a given space. The optimal solution to the problem is a piecewise affine (PWA) function with a polyhedral partition of the entire space of varying parameters. The optimality of the mp-MILP approach is ensured by exhaustiveness, as in any standard Branch and Bound algorithm.

We observe the mp-MILP approach is suitable for real-time systems that must handle workload changes by switching among discrete rates. The key advantage of the multiparametric programming is that, while the offline compo-

nent may have a high time complexity, the online step can generate optimal solutions efficiently. As a result, the optimal solution can be computed quickly in response to workload changes. This characteristic makes it very suitable for the discrete rate adaptation problem. To our knowledge MPRA is the first instantiation of the general multiparametric programming approach in the area of real-time systems.

### 4.2 Problem Transformation

The key step in the design of MPRA is to transform the discrete rate adaptation problem presented in Section 3.2 to an mp-MILP problem, after which the mp-MILP approach is adopted to solve the problem. In the following, we will first introduce the notations for the problem transformation and then use an example to clarify the definitions.

In order to transform the problem, we first introduce a rate adaptation vector  $X$  with  $\bar{m}$  elements, where  $\bar{m} = \sum_{1 \leq i \leq m} k_i$  and  $k_i$  is the number of non-zero rate choices of *adaptable* task  $T_i$ , to represent the rate configuration of the system such that

$$x_l = \begin{cases} 1 & \text{if } T_i \text{ is configured with } r_i^{(j)} \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

where  $l = \sum_{1 \leq s < i} k_s + j$ ,  $1 \leq i \leq m$ , and  $1 \leq j \leq k_i$ . Each 0-1 element in  $X$  corresponds to one non-zero rate choice of some task in an appropriate order. The task rate vector  $R$  can be obtained by  $R = ZX$ , where  $Z$  is an  $m \times \bar{m}$  matrix such that

$$z_{il} = \begin{cases} r_i^{(j)} & \text{if } \sum_{1 \leq s < i} k_s < l \leq \sum_{1 \leq s \leq i} k_s \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

where  $1 \leq i \leq m$ ,  $1 \leq l \leq \bar{m}$ , and  $j = l - \sum_{1 \leq s < i} k_s$ . Each row in  $Z$  is associated with one *adaptable* task and contains the information of the non-zero rate options for the task.

We then introduce an  $n \times m$  matrix  $H$ , where  $h_{ij} = \sum_{T_j l \in S_{ij}} c_{jl}$ , i.e., the total execution time of task  $T_j$ 's subtasks on processor  $P_i$ , and  $h_{ij} = 0$  if no subtask of  $T_j$  is allocated on processor  $P_i$ . The model that characterizes the relationship between  $U$  and  $X$  is given by  $U = HZX + D$ .

To describe the relationship between  $Q_s$  and  $X$ , we introduce a vector  $\bar{Q}$  such that  $\bar{q}_l = w_i q_i^{(j)}$  where  $l = \sum_{1 \leq s < i} k_s + j$ ,  $1 \leq i \leq m$ , and  $1 \leq j \leq k_i$ . Each element in  $\bar{Q}$  corresponds to one non-zero rate choice of some task. Thus, the system utility is calculated by  $Q_s = \bar{Q}X$ . By denoting  $D_N = B - D$  and  $G = HZ$ , we re-formulate the discrete rate adaptation problem as following:

$$\min_X (-\bar{Q}X) \quad (6)$$

subject to

$$GX \leq D_N \quad (7)$$

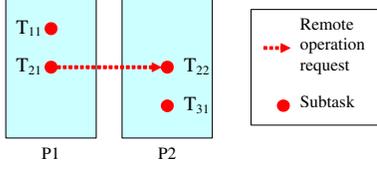


Figure 1. An Example Workload

$$x_i \in \{0, 1\}, 1 \leq i \leq \bar{m} \quad (8)$$

$$\sum_{\sum_{1 \leq s < i} k_s < j \leq \sum_{1 \leq s \leq i} k_s} x_j \leq 1, 1 \leq i \leq m \quad (9)$$

The constraint (7) enforces desired CPU utilization bounds on all processors. The constraint (8) shows that each task supports only a set of task rate choices. For each task only one rate choice can be selected at a time, which is ensured by the constraint (9).

Considering  $D_N$  as the varying parameter vector and  $X$  as the optimization vector, we have transformed the discrete rate adaptation problem to an mp-MILP problem.

**Example:** Suppose there are two processors and three *adaptable* tasks in the system. As shown in Figure 1,  $T_1$  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  $P_2$ . Suppose each task has two non-zero rate options. Then  $\bar{m} = 6$  and  $X = [x_1 \ x_2 \ x_3 \ x_4 \ x_5 \ x_6]^T$ . After the transformation, we have  $H = \begin{bmatrix} c_{11} & c_{21} & 0 \\ 0 & c_{22} & c_{31} \end{bmatrix}$ ,

$$Z = \begin{bmatrix} r_1^{(1)} & r_1^{(2)} & 0 & 0 & 0 & 0 \\ 0 & 0 & r_2^{(1)} & r_2^{(2)} & 0 & 0 \\ 0 & 0 & 0 & 0 & r_3^{(1)} & r_3^{(2)} \end{bmatrix},$$

$$\bar{Q} = [w_1 q_1^{(1)} \ w_1 q_1^{(2)} \ w_2 q_2^{(1)} \ w_2 q_2^{(2)} \ w_3 q_3^{(1)} \ w_3 q_3^{(2)}].$$

The constraint (9) can be described by

$$\begin{bmatrix} 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 \end{bmatrix} X \leq \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}.$$

### 4.3 Design of MPRA

After transforming the discrete rate adaptation problem to an mp-MILP problem, we present the MPRA algorithm that can produce *optimal* rate adaptation solutions online in response to workload changes such as dynamic arrival and departure of *unadaptable* tasks. As shown in Figure 2, MPRA has both offline part and online part. In the following, we present the functionality of each component in detail.

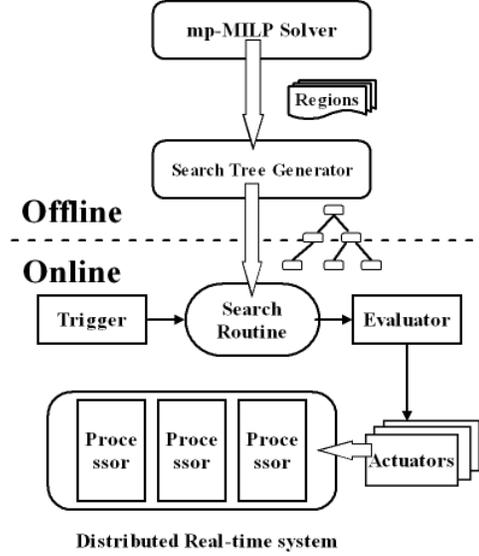


Figure 2. Overview of MPRA

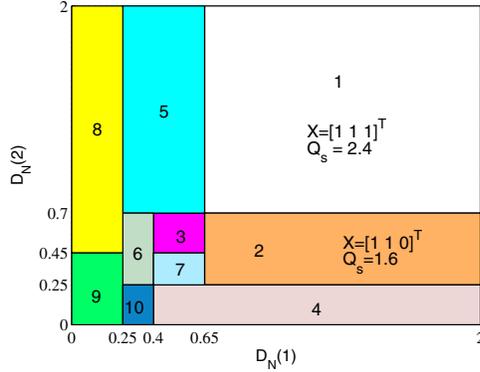
#### 4.3.1 Offline Components

The offline part of MPRA including an mp-MILP Solver and a Search Tree Generator only executes once before the system starts running. It first invokes the mp-MILP Solver to generate the optimal solution, a piecewise affine (PWA) function, and then calls the Search Tree Generator to build a binary tree for the representation of the PWA function.

**mp-MILP Solver:** MPRA invokes the mp-MILP Solver to divide the  $n$ -dimensional space of  $D_N$  into multiple regions and generates the PWA function which expresses  $X$  as a linear function of  $D_N$  for each region. The mp-MILP Solver implements a Branch and Bound algorithm that recursively fixes the 0-1 variables in  $X$  and builds an enumeration tree to generate the PWA function. Each node in the tree corresponds to an intermediate mp-MILP problem with all remaining 0-1 variables. The space of  $D_N$  to be considered for this intermediate problem is defined as the set of regions found for the parent node. At each node, an mpLP problem is solved by relaxing the 0-1 variables as continuous variables in  $[0,1]$ . The solution of a non-leaf node is a lower bound of any integer solution to the intermediate mp-MILP problem. The solution of a leaf node, where all 0-1 variables have been fixed, is an *integer* solution of the final mp-MILP problem in a set of regions. At any level of the tree, the current solution is compared with the upper bound to eliminate parts of the space of  $D_N$  defined for the remaining nodes. Note that the integer solution at each leaf node is feasible (i.e., meets the utilization constraints), but may not be optimal for the final mp-MILP problem in terms of system utility, because the regions for different leaf nodes can overlap with each other. This is undesirable because, for some given value of  $D_N$  that belongs to the intersection

**Table 1. Parameters in the Example Workload**

$T_{ij}$	$c_{ij}$ (ms)	$r_i^{(1)}$ (1/ms)	$q_i^{(1)}$	$\omega_i$
$T_{11}$	20	1/50	0.6	1
$T_{21}$	20	1/80	1.0	1
$T_{22}$	20			
$T_{31}$	45	1/100	0.8	1



**Figure 3. Partition of the Parameter Space with the Example Workload**

of multiple regions, the online part would have to compare the solutions in all those regions to find the optimal one. To facilitate efficient online calculation, the Solver removes the overlap among the regions for all leaf nodes by dividing them into non-overlapping subregions each corresponding to the optimal solution.

**Search Tree Generator:** It generates a binary tree data structure for the representation of the PWA function generated by the mp-MILP Solver. Each node in the tree corresponds to a polyhedron which consists of a set of non-overlapping regions generated by the mp-MILP Solver. An intermediate node contains the affine function for one selected hyperplane that is best for balancing the node’s left and right child in terms of the number of linear functions. Each leaf node maintains one unique linear function that can be evaluated to obtain the optimal solution for any given value of  $D_N$  that belongs to the polyhedron corresponding to this node. For a given  $D_N$  the online part only evaluates one linear inequality at each level and then select the left or right sub-tree to continue based on the sign. With the help of the binary tree, the time of the online evaluation of the PWA function becomes *logarithmic* in the number of regions.

We implemented the offline part of MPRA using the MPT toolbox [16], which provides an mp-MILP solver [4] and a binary tree generator [33].

**Example:** Recall the example workload shown in Figure 1 to demonstrate how the offline part of MPRA works.

Suppose each task only has one non-zero rate choice. Task parameters are given in Table 1. The space of the varying parameters is a 2-dimensional box,  $0 \leq D_N(1) \leq 2$  and  $0 \leq D_N(2) \leq 2$ , for this concrete problem. The mp-MILP Solver divides the entire space into 10 regions (see Figure 3). Each region corresponds to one optimal integer solution. For example,  $X = [1 \ 1 \ 0]^T$  for region 2, i.e., for any given value of  $D_N$  in region 2, only  $T_1$  and  $T_2$  will be admitted in order to maximize the utility while meeting the utilization constraints. The Search Tree Generator generates a binary tree with depth equal to 4 to represent the optimal solution. With the help of the binary tree, we only need to evaluate at most 4 linear inequalities to locate the region that contains a given value of  $D_N$  and obtain the optimal rate configuration.

### 4.3.2 Online Components

Online rate adaptation is triggered by dynamic arrival and departure of *unadaptable* tasks, such as mission critical tasks with fixed rates. Online rate adaptation works as following:

1. **Trigger:** The Trigger calculates  $D$  based on the execution times and rates of the arrived tasks or departed tasks and sends the new value of  $D$  to the Search Routine.
2. **Search Routine:** After receiving  $D$  from the Trigger, the Search Routine traverses the binary tree to locate the region that the current value of  $D_N$  belongs to, and then passes the region number to the Evaluator.
3. **Evaluator:** The Evaluator computes the new value of  $X$  by evaluating the linear function of the region located by the Search Routine. It then sends the new value of  $X$  to Actuators.
4. **Actuator:** the Actuators change the task rates based on the new value of  $X$ . If the new task rate of  $T_i$  is zero,  $T_i$  will be evicted.

**Example:** For the example used in Section 4.3.1, suppose all 3 tasks are running and the current value of  $D_N$  is  $[0.7 \ 0.5]^T$ . The Search Routine goes through the binary tree and ends with region 2. The evaluator then obtains the optimal integer solution  $X = [1 \ 1 \ 0]^T$  corresponding to region 2 and passes it to the actuators on  $P_1$  and  $P_2$ .  $T_3$  will be evicted by the actuator on  $P_2$  according to the new value of  $X$ .

### 4.4 Complexity Analysis

In this section we analyze the complexity of the MPRA algorithm. The complexity of the offline part is exponential in the number of decision variables [25], which is equal

to  $\bar{m}$  for discrete rate adaptation, where  $\bar{m} = \sum_{1 \leq i \leq m} k_i$ ,  $m$  is the number of the *adaptable* tasks, and  $k_i$  is the number of non-zero rates of task  $T_i$ . Note that the exponential complexity is unavoidable in order to get optimal solutions due to the fact that the discrete rate adaptation is an NP-hard problem. A key advantage of MPRA is that it only incurs exponential complexity in the *offline* part which is not time critical and can use significant computing resources. In the following, our analysis focuses on the online search routine and the evaluation of the explicit solution, which dominate the online complexity of MPRA.

The complexity of the online search routine depends on  $N_r$ , the number of non-overlapping regions generated by the mp-MILP Solver. We first analyze the mp-MILP algorithm to calculate  $N_r$ . The mp-MILP Solver implements the Branch and Bound algorithm presented in [4]. There will be  $2^{\bar{m}}$  leaf nodes in the enumeration tree. For each leaf node, all  $\bar{m}$  binary variables have been fixed and the problem is relaxed to an mpLP problem. Based on the results in [5], the upper bound to the number of regions for one leaf node is  $n_r \leq n + 1$ , where  $n$  is the number of processors.

The optimal PWA function of the mp-MILP problem is obtained by removing the overlap among the regions for all leaf nodes. One such region can be divided into at most  $2^{\bar{m}}$  non-overlapping regions because it can be associated with at most  $2^{\bar{m}}$  solutions. After eliminating the intersection among different regions, we get all  $N_r$  non-overlapping regions, which represent a partition of the entire space of  $D_N$ .  $N_r$  is bounded by  $N_r \leq 2^{\bar{m}} \times n_r \times 2^{\bar{m}} \leq (n+1)2^{2\bar{m}}$ .

The binary tree generated by the Search Tree Generator reduces the complexity of online region search. For a given  $D_N$  we only evaluate one linear inequality at each level, which incurs  $n$  multiplications,  $n$  additions and 1 comparison. Traversing the tree from the root to the bottom, we will end up with a leaf node that gives us the optimal solution. Then we need  $2\bar{m}n$  arithmetic operations for the explicit solution evaluation. According to the result in [33], the depth of the binary tree,  $d$ , is given by  $d = \lceil \frac{\ln N_r}{\ln 1/\alpha} \rceil \leq \lceil \frac{2\bar{m} \ln 2 + \ln(n+1)}{\ln 1/\alpha} \rceil$ , where  $0.5 \leq \alpha < 1$ . The constant  $\alpha$  is related to how inbalance the binary tree is. A conservative estimate of  $\alpha$  is 2/3 based on the result in [33]. So the worst-case number of arithmetic operations required for online search and evaluation is  $(2n+1)d + 2\bar{m}n$ . Let  $k = \max\{k_1, \dots, k_m\}$ . Then  $\bar{m} = \sum_{1 \leq i \leq m} k_i \leq km$ . Thus, MPRA has time complexity  $O(n \log(n)) + O(mn)$ , where  $m$  is the number of *adaptable* tasks and  $n$  is the number of the processors.

## 5 Experimentation

In this section, we present simulation results for discrete rate adaptation. Our simulation environment is composed of an event-driven simulator implemented in C++ and the

online part of MPRA. The offline pre-processing of MPRA is done in MATLAB.

In our simulation, the subtasks on each processor are scheduled by the Rate Monotonic scheduling (RMS) algorithm [22]. Each task's end-to-end deadline is  $m_i/r_i$ , where  $m_i$  is the number of subtasks of task  $T_i$  and  $r_i$  is the current rate of the task. The deadline of each task is evenly divided into subdeadlines for its subtasks. The resultant subdeadline of each subtask  $T_{ij}$  equals to its period,  $1/r_i$ . Hence we choose the schedulable utilization bound of RMS [22] as the utilization bound on each processor:  $b_i = n_i(2^{1/n_i} - 1)$ ,  $1 \leq i \leq n$ , where  $n_i$  is the number of subtasks on  $P_i$ . MPRA can also be used with other scheduling policies and their suitable utilization bounds.

We develop a workload generator to create end-to-end tasks and the workload for each set of the experiments. In our simulation, each *adaptable* task has three rate options. We assume every *adaptable* task can be evicted, i.e.,  $r_i^{(0)} = 0$ ,  $1 \leq i \leq m$ .  $r_i^{(1)}$  of task  $T_i$  is the reciprocal of task period  $\tau_i$ , which follows a uniform distribution between 100 ms and 1100 ms. The ratio  $r_i^{(2)}/r_i^{(1)}$  is uniformly distributed between 1.5 and 3. We intentionally choose a small set of rate choices for each task in order to stress-test MPRA's capability to support *discrete* rate options. The fewer rates per task, the more significantly does the problem deviate from the continuous case. The task utility value  $q_i^{(0)}$  of  $T_i$  when the task is evicted is zero and  $q_i^{(1)}$  at rate  $r_i^{(1)}$  is randomly generated using a uniform distribution between 0.5 and 2. The ratio of the utilities at different rates,  $q_i^{(2)}/q_i^{(1)}$  is uniformly distributed between 1.5 and 3. All weights are set to 1 for simplicity in our simulation, i.e.,  $w_i = 1$ ,  $1 \leq i \leq m$ . The number of subtasks of each task ranges from 1 to 4 and all subtasks are randomly allocated on all processors. The worst-case execution time  $c_{ij}$  of subtask  $T_{ij}$  is obtained by  $c_{ij} = u_{ij}\tau_i$ , where  $u_{ij}$ , the utilization of  $T_{ij}$ , is uniformly distributed from 0.05 to 0.2.

### 5.1 Baselines and Performance Metrics

We compare MPRA against three existing algorithms: **bintprog**, **amrmd1** [18], and **amrmd\_dp** [12]. **bintprog** is a binary integer linear programming solver provided by the commercial Optimization Toolbox of MATLAB 7. **bintprog** is a representative optimization solver that can produce optimal solutions, which is used to validate the optimality of MPRA. **amrmd1** and **amrmd\_dp**, where **amrmd** stands for Approximate Multi-Resource Multi-Dimensional Algorithm, are two representative efficient heuristic algorithms for utility optimization in real-time systems. **amrmd\_dp** can perform better than **amrmd1** in terms of utility at the cost of longer execution time than **amrmd1**.<sup>2</sup> How-

<sup>2</sup>The authors also present another algorithm called **amrmd\_cm** to address the co-located point problem of **amrmd1** in [12]. It performs exactly

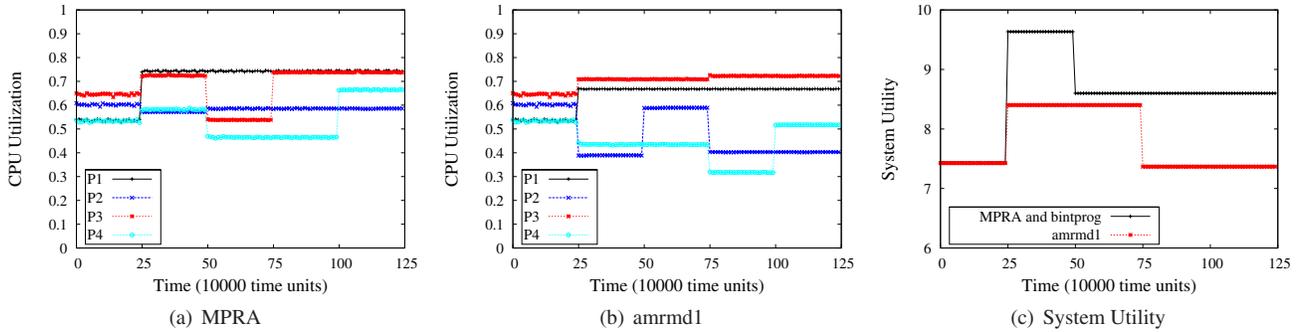


Figure 4. System Performance under Separate Task Arrivals

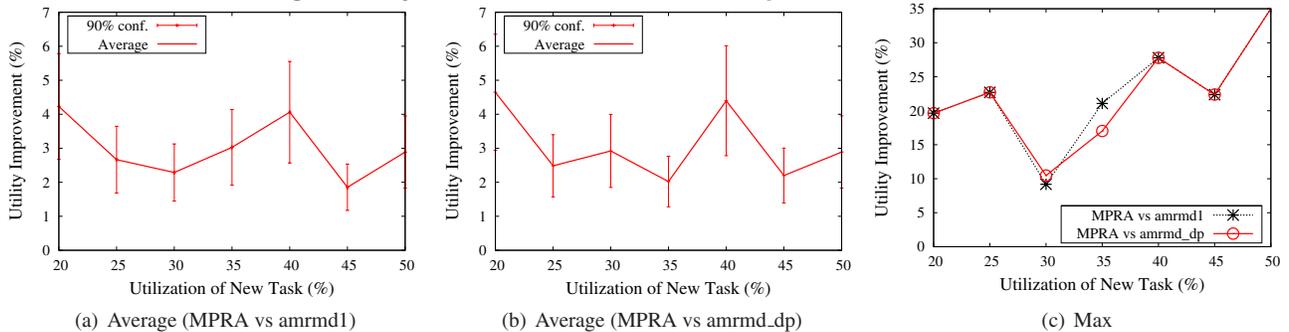


Figure 5. Utility Improvement over Heuristics

ever, **amrmd1** and **amrmd\_dp** may produce sub-optimal solutions and do not have theoretical error bounds as mentioned in [19].

In our experiments, online adaptation operations are triggered by arrivals of *unadaptable* tasks. The performance metric used throughout the simulation is *utility improvement*,  $\delta$ , which is defined by  $\delta = (Q_{MPRA} - Q_b)/Q_b$ , where  $Q_{MPRA}$  and  $Q_b$  are the system utilities produced by MPRA and a baseline algorithm, respectively, after they make the online adaptation in response to the same new task arrivals.

In order to evaluate the efficiency of MPRA, we also investigate its online execution time and compare it with three baselines. The execution times are measured on a 2.52GHz Pentium 4 PC with 1 GB RAM. To achieve fine grained measurements, we use the high resolution timer *gethrtime* provided by ACE [9]. This function uses an OS-specific high-resolution timer that returns the number of clock cycles since the CPU is powered up or reset. The *gethrtime* function has a low overhead and is based on a 64 bit clock cycle counter on Pentium processors. To estimate the average computation overhead of an online adaptation operation, we run each online execution for 100 times as a subroutine. The result is then divided by 100 to get the execution time of a single execution.

the same as **amrmd1** here because no co-located points exist in the discrete rate adaptation problem.

## 5.2 Simulation Results

In the simulation each workload includes 6 end-to-end *adaptable* tasks executing on 4 processors. The results are based on 20 randomly generated workloads. In the following, we present two sets of experiments to evaluate the performance of the four algorithms in the presence of arrivals of *unadaptable* tasks. In our experiments, each *adaptable* task has two non-zero rate choices. All tasks in the workload are running at the lower rate in the beginning. The new arrival tasks are *unadaptable* tasks that must be executed at the cost of other tasks.

In the first set of experiments, an *unadaptable* task arrives at each processor at 250000, 500000, 750000, and 1000000 time unit, respectively, which triggers online rate adaptation four times. The CPU utilization of each new task is 0.2. As shown in Figure 4, all algorithms maintain acceptable utilizations on all processors in face of new task arrivals. One important observation is that MPRA and **bintprog** always produce the same optimal rates and hence achieve the same system utility in all the experiments. These results are consistent with the optimality of MPRA. Both MPRA and **bintprog** generate optimal rates that result in higher system utilities than **amrmd1** in response to new task arrivals (see Figure 4(c)). For instance, MPRA has a better achievement in term of utility than **amrmd1** does (9.63 versus 8.41) when both algorithms try to fully utilize the system resource and maximize the system utility

by changing task rates at 250000 time unit. Another example is, after the new task arrives on  $P_3$  at 750000 time unit, the utility achieved by MPRA remains unchanged while the utility achieved by **amrmd1** decreases significantly. As shown in Figure 4(a), the utilizations on other three processors are not affected by the new task arrival on  $P_3$  with the help of MPRA, i.e., MPRA admits the new task without degrading the system performance. In contrast, the utilizations on  $P_1$  and  $P_4$  decrease after 750000 time units in Figure 4(b), which means **amrmd1** has to reduce the rates of *adaptable* tasks in order to handle the task arrival on  $P_3$ .

We run the other set of experiments by varying the CPU utilization of the new arrival task from 0.2 to 0.5. Four identical new tasks are activated after 250000 time units on four processors simultaneously. Consequently, online rate adaptation is triggered to maximize system utility while enforcing the utilization bounds. We repeated the same set of experiments for all 20 workloads. Figure 5 plots the average and maximum *utility improvements* achieved by MPRA over **amrmd1** and **amrmd\_dp** as the utilization of the new *unadaptable* task increases from 0.2 to 0.5. Each data point and the corresponding confidence interval in Figure 5 are calculated from all 20 *utility improvement* results obtained in 20 runs under a given workload variation. As shown in Figure 5(a) and 5(b), MPRA consistently achieves same utilities as **bintprog** and outperforms both **amrmd1** and **amrmd\_dp** in terms of system utility. Moreover, MPRA achieves as high as 35% utility improvement over both **amrmd1** and **amrmd\_dp**. Our results demonstrate that, while state-of-the-art heuristics such as **amrmd1** and **amrmd\_dp** may achieve good (but suboptimal) performance on average, they may result in significantly lower system utility in certain cases. This observation is consistent with the fact that the heuristics do not have analytical bounds on the distance from optimal solutions. In contrast, a fundamental benefit of MPRA is that it can always achieve *optimal* system utility in face of workload variations. The analytical guarantee on optimal system utilities can be highly desirable to dynamic mission-critical applications.

The average execution-times of all four approaches are shown in Figure 6. MPRA's online overhead is more than two orders of magnitude lower than that of **bintprog** while generating the same optimal solutions. MPRA remains comparable to **amrmd1** and **amrmd\_dp** in terms of online overhead. For instance, when the new task has an utilization of 30%, MPRA incurs an overhead of about 100 *microseconds*, while **bintprog** needs about 100 *milliseconds* to generate the same optimal rate assignments. The results show that MPRA can provide optimal rate adaptation for end-to-end tasks with comparable online overhead as efficient suboptimal heuristics. Note that the overhead introduced by MPRA is about 100 *microseconds*, which is negligible when compared to both (i) the time scale of dead-

lines and periods of end-to-end tasks in many distributed real-time systems and (ii) the overhead of more than 20 *milliseconds* incurred by adjusting task rates on middleware systems [34].

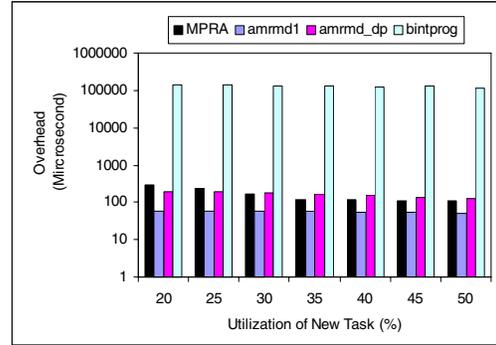


Figure 6. Online Overhead

## 6 Conclusions and Future Work

We have developed the MPRA algorithm for optimal and efficient discrete rate adaptation in distributed real-time systems. We first transform the discrete rate adaptation problem as an mp-MILP problem. We then present the design and complexity analysis of MPRA which proves that MPRA can reduce its online complexity to *polynomial* time through offline preprocessing. Simulation results demonstrate that MPRA maximizes the system utility in face of workload variations, with the online execution time more than two orders of magnitude lower than a representative optimization solver. Moreover, it consistently outperforms efficient heuristics in terms of system utility at comparable online overhead. While we focus on discrete rate adaptation in this paper, the multiparametric approach may be applicable to a broad range of adaptive systems with discrete configurations. In the future we plan to extend our work to other online adaptation mechanisms such as task reallocation or dynamic voltage scaling.

In this paper we evaluate the performance of MPRA in response to dynamic arrivals of tasks with fixed rates. Our next step is to explore other workload changes such as execution time variations. In the current implementation MPRA deals with workload changes that can be calculated explicitly. Our approach may be combined with event-driven feedback control to deal with uncertainties in system workload based on measured CPU utilization. The extension is part of our future work.

## References

- [1] T. Abdelzaher, J. Stankovic, C. Lu, R. Zhang, and Y. Lu. Feedback performance control in software services. *IEEE*

- Control Systems*, 23(3), June 2003.
- [2] T. F. Abdelzaher, E. M. Atkins, and K. G. Shin. QoS Negotiation in Real-Time Systems and Its Application to Automated Flight Control. *IEEE Transactions on Computers*, 49(11):1170–1183, 2000.
  - [3] T. F. Abdelzaher, K. G. Shin, and N. Bhatti. Performance guarantees for Web server end-systems: A control-theoretical approach. *IEEE Transactions on Parallel and Distributed Systems*, 13(1):80–96, 2002.
  - [4] J. Acevedo and E. Pistikopoulos. An Algorithm for Multiparametric Mixed-integer Linear Programming Problems. *Operations Research Letters*, 24:139–148(10), 1999.
  - [5] A. Bemporad, F. Borrelli, and M. Morari. Model Predictive Control Based on Linear Programming - The Explicit Solution. *IEEE Transactions on Automatic Control*, 47(12):1974–1985, Dec. 2002.
  - [6] S. Brandt, G. Nutt, T. Berk, and J. Mankovich. A Dynamic Quality of Service Middleware Agent for Mediating Application Resource Usage. In *IEEE Real-Time Systems Symposium*, page 307, Washington, DC, USA, 1998.
  - [7] S. A. Brandt and G. J. Nutt. Flexible Soft Real-Time Processing in Middleware. *Real-Time Systems*, 22(1-2):77–118, 2002.
  - [8] G. C. Buttazzo, G. Lipari, M. Caccamo, and L. Abeni. Elastic Scheduling for Flexible Workload Management. *IEEE Trans. Comput.*, 51(3):289–302, 2002.
  - [9] Center for Distributed Object Computing. The ADAPTIVE Communication Environment (ACE). [www.cs.wustl.edu/~schmidt/ACE.html](http://www.cs.wustl.edu/~schmidt/ACE.html), Washington University.
  - [10] A. Cervin, J. Eker, B. Bernhardsson, and K.-E. Årzén. Feedback-Feedforward Scheduling of Control Tasks. *Real-Time Systems*, 23(1-2):25–53, 2002.
  - [11] T. Gal and J. Nedoma. Multiparametric Linear Programming. *Management Science*, 18:406–442, 1972.
  - [12] S. Ghosh, R. R. Rajkumar, J. Hansen, and J. Lehoczky. Scalable Resource Allocation for Multi-Processor QoS Optimization. In *Proceedings of the 23rd International Conference on Distributed Computing Systems*, pages 174–183, Washington, DC, USA, 2003.
  - [13] O. H. Ibarra and C. E. Kim. Fast Approximation Algorithms for the Knapsack and Sum of Subset Problems. *Journal of the ACM*, 22(4):463–468, 1975.
  - [14] B. Kao and H. Garcia-Molina. Deadline assignment in a distributed soft real-time system. *IEEE Transactions on Parallel and Distributed Systems*, 8(12):1268–1274, 1997.
  - [15] X. Koutsoukos, R. Tekumalla, B. Natarajan, and C. Lu. Hybrid Supervisory Utilization Control of Real-Time Systems. In *IEEE Real-Time and Embedded Technology and Applications Symposium*, pages 12–21, 2005.
  - [16] M. Kvasnica, P. Grieder, and M. Baotić. Multi-Parametric Toolbox (MPT), 2004.
  - [17] C. Lee, J. Lehoczky, R. Rajkumar, and D. P. Siewiorek. On Quality of Service Optimization with Discrete QoS Options. In *IEEE Real Time Technology and Applications Symposium*, pages 276–286, 1999.
  - [18] C. Lee, J. P. Lehoczky, D. P. Siewiorek, R. Rajkumar, and J. P. Hansen. A Scalable Solution to the Multi-Resource QoS Problem. In *IEEE Real-Time Systems Symposium*, pages 315–326, 1999.
  - [19] C. Lee and D. Siewiorek. An Approach for Quality of Service Management. Technical Report CMU-CS-98-165, Computer Science Department, CMU, 1998.
  - [20] C.-G. Lee, C.-S. Shih, and L. Sha. Online QoS Optimization Using Service Classes in Surveillance Radar Systems. *Real-Time Systems*, 28(1):5–37, 2004.
  - [21] J. P. Lehoczky. Fixed Priority Scheduling of Periodic Task Sets with Arbitrary Deadlines. In *IEEE Real-Time Systems Symposium*, pages 201–213, 1990.
  - [22] C. Liu and J. Layland. Scheduling Algorithms for Multiprogramming in a Hard-Real-Time Environment. *Journal of the ACM*, 20(1):46–61, 1973.
  - [23] C. Lu, X. Wang, and X. Koutsoukos. Feedback Utilization Control in Distributed Real-Time Systems with End-to-End Tasks. *IEEE Transactions on Parallel Distributed Systems*, 16(6):550–561, June 2005.
  - [24] S. Martello and P. Toth. *Knapsack problems: algorithms and computer implementations*. John Wiley & Sons, Inc., New York, NY, USA, 1990.
  - [25] K. G. Murty. Computational Complexity of Parametric Linear Programming. *Mathematical Programming*, 19(1):213–219, December 1980.
  - [26] M. D. Natale and J. Stankovic. Dynamic End-to-end Guarantees in Distributed Real-time Systems. In *IEEE Real-Time Systems Symposium*, 1994.
  - [27] R. Rajkumar, C. Lee, J. Lehoczky, and D. Siewiorek. A Resource Allocation Model for QoS Management. In *IEEE Real-Time Systems Symposium*, pages 298–307, Dec. 1997.
  - [28] R. Rajkumar, C. Lee, J. P. Lehoczky, and D. P. Siewiorek. Practical Solutions for QoS-Based Resource Allocation. In *IEEE Real-Time Systems Symposium*, pages 296–306, 1998.
  - [29] S. Sahni. Approximate Algorithms for the 0/1 Knapsack Problem. *J. ACM*, 22(1):115–124, 1975.
  - [30] D. C. Steere, A. Goel, J. Gruenberg, D. McNamee, C. Pu, and J. Walpole. A Feedback-driven Proportion Allocator for Real-Rate Scheduling. In *Operating Systems Design and Implementation*, pages 145–158, 1999.
  - [31] J. Sun and J. W.-S. Liu. Synchronization Protocols in Distributed Real-Time Systems. In *International Conference on Distributed Computing Systems*, pages 38–45, 1996.
  - [32] H. Tokuda and T. Kitayama. Dynamic QoS Control based on Real-Time Threads. In *NOSSDAV 93*, pages 114–123, London, UK, 1994. Springer-Verlag.
  - [33] P. Tondel, T. A. Johansen, and A. Bemporad. Evaluation of Piecewise Affine Control via Binary Search Tree. *Automatica*, 39:945–950, 2003.
  - [34] 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*, 80(7):938–950, 2007.
  - [35] X. Wang, D. Jia, C. Lu, and X. Koutsoukos. Decentralized Utilization Control in Distributed Real-Time Systems. In *IEEE Real-Time Systems Symposium*, pages 133–142, 2005.