

Recovering from Overload in Multicore Mixed-Criticality Systems

Jeremy P. Erickson, Namhoon Kim, and James H. Anderson

Department of Computer Science, University of North Carolina at Chapel Hill*

Abstract

The MC^2 mixed-criticality framework has been previously proposed for mixing safety-critical hard real-time (HRT) and mission-critical soft real-time (SRT) software on the same multicore computer. This paper focuses on the execution of SRT software within this framework. When determining SRT guarantees, jobs are provisioned based on a provisioned worst-case execution time (PWCET) that is not very pessimistic and could be overrun. In this paper, we propose a mechanism to recover from the overload created by such overruns. Specifically, we propose a modification to the previously proposed G-EDF-like (GEL) class of schedulers that uses virtual time to increase task periods. We also experimentally determine how long it takes to return to normal behavior after a transient overload.

1 Introduction

Future cyber-physical systems will require mixing tasks of varying importance. For example, future unmanned aerial vehicles (UAVs) will require more stringent timing requirements for adjusting flight surfaces than for long-term decision-making [11]. The *mixed-criticality* (MC) framework MC^2 has been previously proposed in order to allow workloads of differing criticalities to be simultaneously supported on a single multicore machine [11, 17]. Using a single machine allows reductions in size, weight, and power.

In any MC system, there are a number of *criticality levels*. For example, MC^2 has four criticality levels, denoted A (highest) through D (lowest). Each MC task is assigned a distinct criticality level. When analyzing an MC system, each task is assigned a separate *provisioned worst-case execution time* (PWCET) for each criticality level. Guarantees are provided for level- ℓ tasks by assuming that no task with criticality at or above level ℓ exceeds its level- ℓ PWCET. For example, when analyzing level C, level-A, -B, and -C tasks are considered using level-C PWCETs.

As noted by Burns and Davis [6], most proposed MC frameworks do not provide any guarantees for a given level ℓ if any job exceeds its level- ℓ PWCET. This assumption could be highly problematic in practice. For example, suppose that a level-A flight-control job on a UAV exceeds its level-C PWCET. It would be very undesirable if no guarantees could be provided for level-C mission control tasks

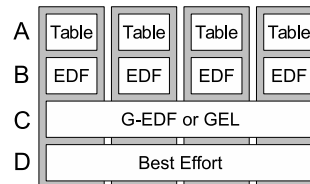


Figure 1: MC^2 architecture.

from that point forward.¹ The purpose of this paper is to provide guarantees in MC^2 in such situations.

Specifically, we consider response-time behavior for tasks at level C in MC^2 . The architecture of MC^2 as described by Herman et al. [11] is depicted in Fig. 1: levels A and B are scheduled on a per-processor basis using table-driven and EDF scheduling, respectively. Level C was proposed by Mollison et al. [17] to be scheduled using the global earliest-deadline-first (G-EDF) scheduler, which provides *bounded response times* but may not meet all deadlines. Erickson et al. [9] demonstrated that the more general class of *G-EDF-like* (GEL) schedulers [15] can yield better response-time bounds. Therefore, we consider general GEL schedulers here.

Contributions. When any job at or above level C overruns its level-C PWCET, the system at level C may be *overloaded*, compromising level-C guarantees. Using the MC^2 framework, a task may have its per-job response times permanently increased as a result of even a single overload event, and multiple overload events could cause such increases to build up over time. Examples of conditions that could cause this to happen are presented in Sec. 2. As a result, we must alter scheduling decisions to attempt to recover from transient overload conditions. In this paper, we propose a scheme that does so by scaling task inter-release times and modifying scheduling priorities. We further present an implementation of this scheme, including both in-kernel and userspace components. We also provide experimental results based on an actual implementation that demonstrates that this scheme can effectively recover from overload.

Comparison to Related Work. Other techniques for managing overload have been provided in other settings, although most previously proposed techniques either focus exclusively on uniprocessors [2, 3, 7, 13, 16] or only pro-

*Work supported by NSF grants CNS 1016954, CNS 1115284, CNS 1218693, and CNS 1239135.

¹At the industry session of RTAS 2014, several industry practitioners noted this as a practical concern that had not been adequately addressed in the literature on MC systems.

vide heuristics without theoretical guarantees [10].

Our paper uses the idea of “virtual time” from Zhang [23] (as also used by Stoica et al. [20]), where job separation times are determined using a virtual clock that changes speeds with respect to the actual clock. In our work, we recover from overload by slowing down virtual time, effectively reducing the frequency of job releases. Unlike in [20], we never speed up virtual time relative to the normal underloaded system, so we avoid problems that have previously prevented virtual time from being used on a multiprocessor. To our knowledge, this work is the first to use virtual time in multiprocessor scheduling.

Some past work on recovering from PWCET overruns in MC systems has used techniques similar to ours, albeit in the context of trying to meet all deadlines [12, 18, 19, 21, 22]. Our scheme is also similar to reweighting techniques that modify task parameters such as periods. A detailed survey of several such techniques is provided by Block [4].

Organization. In Sec. 2, we describe the task model used in prior work and show why overload can cause guarantees to be permanently violated. In Sec. 3, we describe our modified task model and scheduler, and discuss how it can be used to recover from overload. In Sec. 4, we describe our implementation, and in Sec. 5, we provide experimental evidence that our scheme is effective.

2 Original MC² And Overload

In this paper, we assume that time is continuous and we consider only the system at level C. In other words, we consider level-A and -B tasks as CPU supply that is unavailable to level C, rather than as explicit tasks. We consider a system $\tau = \{\tau_0, \tau_1, \dots, \tau_{n-1}\}$ of n level-C tasks running on m processors $P = \{P_0, P_1, \dots, P_{m-1}\}$. Each τ_i is composed of a (potentially infinite) series of jobs $\{\tau_{i,0}, \tau_{i,1}, \dots\}$. The release time of $\tau_{i,k}$ is denoted as $r_{i,k}$. We assume that $\min_{\tau_i \in \tau} r_{i,0} = 0$. Each $\tau_{i,k}$ is prioritized on the basis of a *priority point (PP)*, denoted $y_{i,k}$. The time when $\tau_{i,k}$ actually completes is denoted $t_{i,k}^c$, and its *actual execution time* is denoted $e_{i,k}$. We define the *response time* $R_{i,k}$ of $\tau_{i,k}$ as $t_{i,k}^c - r_{i,k}$. We define a job $\tau_{i,k}$ as *pending* at time t if $r_{i,k} \leq t < t_{i,k}^c$.

Under GEL scheduling and the conventional sporadic task model, each task is characterized by a per-job worst-case execution time (WCET) $C_i > 0$, a minimum separation $T_i > 0$ between releases, and a relative PP $Y_i \geq 0$. Using the above notation, the system is subject to the following constraints for every $\tau_{i,k}$:

$$e_{i,k} \leq C_i, \quad (1)$$

$$r_{i,k+1} \geq r_{i,k} + T_i, \quad (2)$$

$$y_{i,k} = r_{i,k} + Y_i. \quad (3)$$

In our work, we consider *provisioned* WCETs due to the mixed-criticality analysis, as discussed in the introduction.

Prior work [14, 17] shows that bounded response times

can be achieved for level-C tasks assuming certain constraints on system-wide and per-task utilizations. To illustrate this property, we depict in Fig. 2(a) a system that only has level-A and -C tasks, with one level-A task per CPU. For level-A tasks, we use the notation (T_i, C_i^C, C_i^A) , where T_i is task τ_i 's period, C_i^C is its level-C PWCET, and C_i^A is its level-A PWCET. For level-C tasks, we use the notation (T_i, Y_i, C_i) . Observe that in Fig. 2(a), no job runs for longer than its level-C PWCET. Under this condition, response times can be bounded using techniques from prior work [14, 17]. In this paper, we typically concern ourselves with response times relative to a job's PP. Under the model we are defining in this section, such a response time can be converted to an absolute response time by adding Y_i . Observe that in Fig. 2(a) some jobs do complete after their PPs; this is allowed by our model. Similarly, some jobs complete after the release of their respective successor jobs.

The particular example in Fig. 2(a) fully utilizes all processors. In the situation depicted in Fig. 2(b), both level-A tasks released at time 12 run for their full level-A PWCETs. Therefore, from the perspective of level C, an overload occurs.² Because the system is fully utilized, there is no “slack” that allows for recovery from overload, and response times are permanently increased. In a system with large utilization, response times could take significant time to settle back to normal, even if they eventually will.

Another cause of overload is depicted in Fig. 3, where there is only a single level-C task. Observe that in Fig. 3(a) τ_1 executes except when both CPUs are occupied by level-A tasks. Therefore, when the overload occurs at time 12 in Fig. 3(b), τ_1 cannot recover despite the frequent presence of slack on the other CPU. This demonstrates that an overload can cause long-running problems due to a single task's utilization, not merely due to the total utilization of the system.

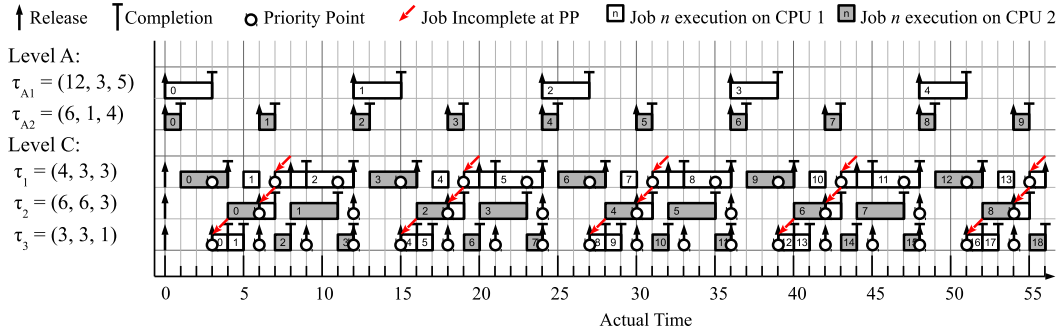
In the next section, we discuss our approach for recovering from overload.

3 Our Modifications

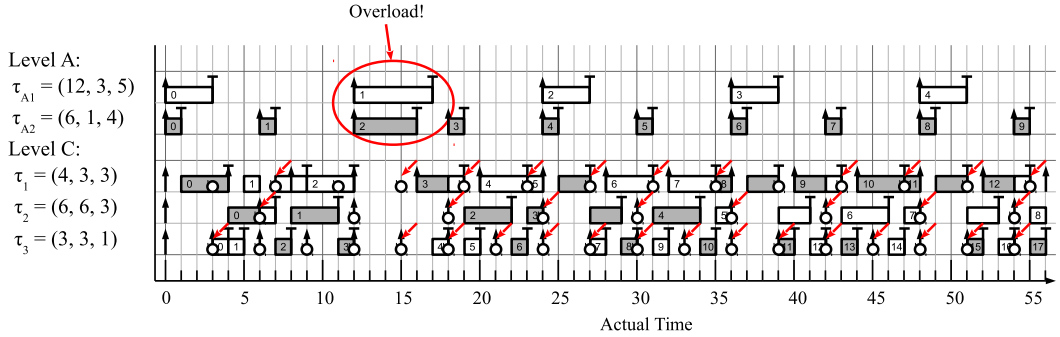
In order to recover from overload, it is necessary to effectively *reduce* task utilizations, to avoid the problems discussed in the previous section. In this paper, we propose to do so by using a notion of *virtual time* (as in [20]), as described in this section.

Our scheme involves a generalized version of GEL scheduling, called *GEL with virtual time (GEL-v)* scheduling, and a generalized version of the sporadic task model, called the *sporadic with virtual time and overload (SVO) model*. Under the SVO model, we no longer assume a particular WCET (thus allowing overload). Therefore, (1) is no

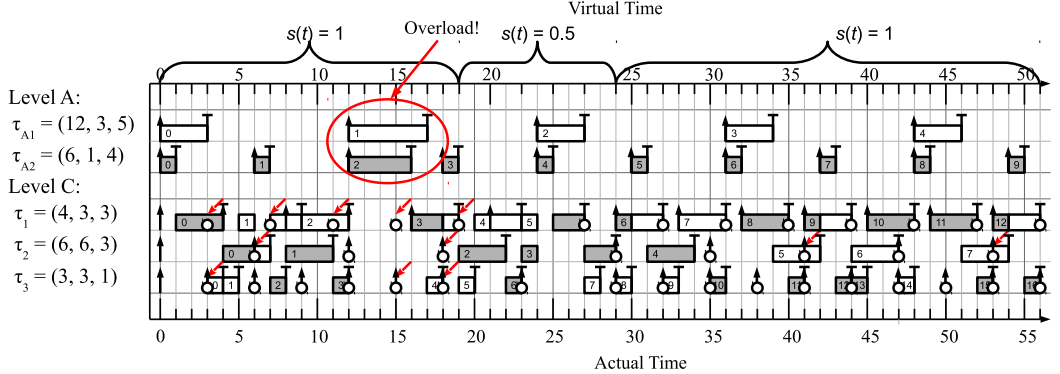
² A similar overload could occur if a level-C task exceeds its level-C PWCET. However, MC² optionally supports the use of execution budgets in order to prevent such an occurrence. While the use of execution budgets would prevent level-A and -B tasks from overrunning their level-A and -B PWCETs, respectively, they can still overrun their level-C PWCETs. Thus, we have chosen examples that provide overload even when execution budgets are used.



(a) Example MC^2 schedule in the absence of overload, illustrating bounded response times.



(b) The same schedule in the presence of overload caused by level-A tasks running for their full level-A PWCETs. Notice that response times of level-C jobs settle into a pattern that is degraded compared to (a). For example, consider $\tau_{2,6}$, which is released at actual time 36. In (a) it completes at actual time 43 for a response time of 7, but in this schedule it does not complete until actual time 46, for a response time of 10.



(c) The same schedule in the presence of overload and the recovery techniques described in Sec. 3. Notice that response times of level-C jobs settle into a pattern that is more like (a) than (b). For example, consider again $\tau_{2,6}$, which is now not released until actual time 41 and completes at actual time 47 for a response time of 6. This is more similar to (a) than to (b).

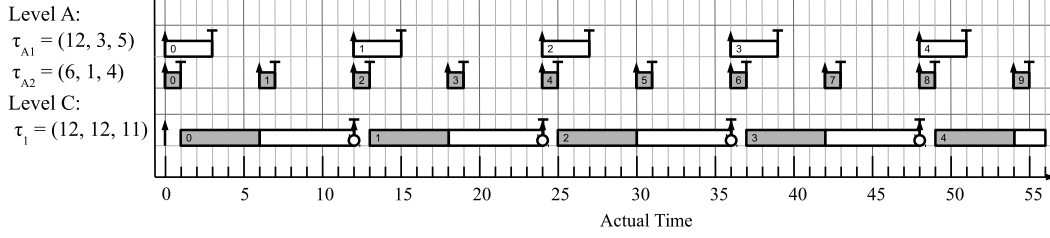
Figure 2: Example MC^2 task system, without and with overload.

longer required to hold.³ Under GEL-v scheduling and the SVO model, we also introduce the use of virtual time, and we define the minimum separation time and relative PP of a task with respect to virtual time after one of its job releases instead of actual time. Virtual time affects *only* level C, not levels A and B. The use of virtual time will allow us to recover from overload. We now introduce our strategy using

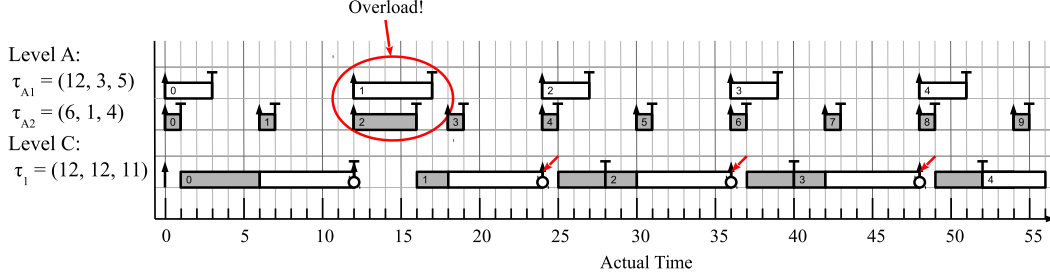
³As mentioned in Footnote 2, execution budgets can be used to restore this assumption at level C, in which case overloads can come only from levels A and B.

the example depicted in Fig. 2(c).

Once an overload occurs, the system can respond by altering virtual time for level C. Virtual time is based on a global speed function $s(t)$. During normal operation of the system, $s(t)$ is always 1. This means that actual time and virtual time progress at the same rate. However, after an overload occurs, the scheduler may choose to select $0 < s(t) < 1$, at which point virtual time progresses more slowly than actual time. In Fig. 2(c), the system chooses to use $s(t) = 0.5$ for $t \in [19, 29]$. As a result, virtual time



(a) Example MC^2 schedule in the absence of overload, illustrating bounded response times.



(b) The same schedule in the presence of overload caused by level-A tasks running for their full level-A PWCEts.

Figure 3: Another example MC^2 task system, without and with overload. See Fig. 2 for key.

progresses more slowly in this interval, and new releases of jobs are delayed. This allows the system to recover from the overload, so at actual time 29, $s(t)$ returns to 1. Observe that job response times are significantly increased after actual time 12 when the overload occurs, but after actual time 29, they are similar to before the overload. In fact, the arrival pattern of level A happens to result in better response times after recovery than before the overload, although this is not guaranteed under a sporadic release pattern.

An actual time t is converted to a virtual time using

$$v(t) \triangleq \int_0^t s(t) dt. \quad (4)$$

For example, in Fig. 2(c), $v(25) = \int_0^{25} s(t) dt = \int_0^{19} 1 dt + \int_{19}^{25} 0.5 dt = 19 + 3 = 22$. Unless otherwise noted, all instants herein (e.g., t , $r_{i,k}$, etc.) are specified in actual time, and all variables *except* T_i and Y_i (defined below) refer to quantities of actual time.

Under the SVO model, (2) generalizes to

$$v(r_{i,k+1}) \geq v(r_{i,k}) + T_i, \quad (5)$$

and under GEL-v scheduling, (3) generalizes to

$$v(y_{i,k}) = v(r_{i,k}) + Y_i. \quad (6)$$

For example, in Fig. 2(c), $\tau_{1,0}$ is released at actual time 0, has its PP three units of (both actual and virtual) time later at actual time 3, and $\tau_{1,1}$ can be released four units of (both actual and virtual) time later at time 4. However, $\tau_{1,5}$ of the same task is released at actual time 21, shortly after the virtual clock slows down. Therefore, its PP is at actual time 27, which is three units of *virtual* time after its release, and the

release of $\tau_{1,6}$ can be no sooner than actual time 29, which is four units of *virtual* time after the release of $\tau_{1,5}$. However, the execution time of $\tau_{1,5}$ is not affected by the slower virtual clock.

In a real system, unlike in our examples so far, level-C jobs will often run for less time than their respective level-C PWCEts. Therefore, it may be unnecessarily pessimistic to initiate overload response whenever a job overruns its level-C PWCEt. Instead, we use the following definition.

Def. 1. τ_i has a nonnegative *response-time tolerance*, denoted ξ_i , relative to each job's PP. A task *meets* its response-time tolerance if $t_{i,k}^c \leq y_{i,k} + \xi_i$, and *misses* it otherwise.

We slow down the virtual clock only after some job misses its response-time tolerance. Ideally, response-time tolerances should be determined based on analytical upper bounds of job response times, in order to guarantee that the virtual clock is never slowed down in the absence of overload. However, for illustration, in Fig. 2(c) we simply use a response-time tolerance of three for each task. Thus, we do not slow down virtual time until some job's completion time is greater than three units of actual time after its PP. At time 18, $\tau_{3,4}$ completes exactly three units after its PP, which is barely within its tolerance, so the virtual clock is not slowed down. However, at time 19, $\tau_{1,3}$ completes four units after its PP, which exceeds the response-time tolerance. Therefore, we slow down the virtual clock at time 19.

We will define *normal behavior* for a system as the situation in which all jobs meet their response-time tolerances. Recall that, as depicted in Fig. 2 above, a system with high utilization may not effectively be able to recover from overload, because there is no slack, and as depicted in Fig. 3 above, a system with a task of high utilization may not be able to effectively recover from overload. As we have

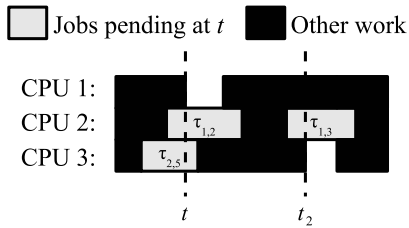


Figure 4: Illustration of “idle normal instant.” If all jobs pending at t meet their response-time tolerances, then t is an idle normal instant. t_2 is referenced in Sec. 4.

just discussed, our technique creates extra slack both in a system-wide sense and in a per-task sense, solving both problems. Therefore, the system eventually returns to normal behavior. We denote the time required to do so as *dissipation time*.

In a technical report [8], we provide theoretical analysis of dissipation time. In that technical report, we first provide analytical upper bounds on response time that can safely be used as response-time tolerances. We also derive an upper bound on dissipation time, called a *dissipation bound*, with respect to these response-time tolerances. As discussed above, our technique causes the system to eventually return to normal behavior, so this bound exists. In this paper, rather than considering theoretical dissipation bounds, we focus on experimentally determining dissipation time at runtime.

In the demand analysis used in our technical report, as in most such analysis, demand is considered beginning at an instant when some processor is idle. If all jobs pending at this time meet their response-time tolerances, then regardless of how that situation arose, the system has returned to normal behavior. Furthermore, the virtual clock can safely be returned to speed 1 after such an instant. Therefore, the system can detect such an instant to determine when to set the virtual-clock speed back to 1. We now define such an instant more formally.

Def. 2. Arbitrary time t is an *idle normal instant* if some processor is idle at t and all jobs pending at t meet their (normal) response-time tolerances.

Our method does not dictate a particular choice of $s(t)$, although in our experiments we consider several such values. Selecting a small value of $s(t)$ will result in a large short-term impact on level-C job releases, but the system will return to normal behavior quickly. Alternatively, selecting a large value of $s(t)$ will result in a lesser short-term impact on job releases, causing only minor delays, but the system will take a longer time in order to return to normal behavior. In our experimental comparison, we quantify these effects and point to proper design decisions.

As suggested by the analysis in our technical report [8], we determine when the system returns to normal behavior by detecting an idle normal instant. Therefore, we return the virtual clock to speed 1 after detecting such a t , which can only be determined when all jobs pending at t are complete. In Fig. 2(c), observe that only CPU 2 is executing work

from actual time 28 to actual time 29. Thus, only τ_2 is pending throughout this interval, or CPU 1 would be executing work. Furthermore, $\tau_{2,4}$ is the only pending job of τ_2 at time 28. Observe that $\tau_{2,4}$ completes at its PP, and thus meets its response-time tolerance of three, at time 29. Therefore, time 28 is an idle normal instant. The system can determine this to be the case at time 29, when $\tau_{2,4}$ completes. Therefore, the virtual clock returns to speed 1 at time 29.

4 Implementation Description

We implemented our scheme by extending the existing MC² implementation that was described in [11]. That implementation is based on LITMUS^{RT} [1], a real-time extension to Linux originally developed at UNC. Source code for our implementation is also available at [1]. Our implementation consists of two components: the scheduler, which is part of the kernel, and a *monitor program*, which runs in userspace. The kernel reports job releases and job completions to the monitor program and provides a system call that the monitor program can use to change the speed of the virtual clock. The speed of the virtual clock does not change between these calls. The kernel is responsible for implementing virtual time, ensuring that the SVO model’s minimum-separation constraints are respected, and making scheduling decisions according to GEL-v scheduling. The monitor program is responsible for determining when virtual-clock speed changes should occur.

Within the kernel, the primary change that we made compared to the prior MC² implementation was the use of virtual time at level C. No changes at levels A or B were required. Pseudocode for the changed functionality is provided in Algo. 1. `now()` is a function that always returns the current actual time.

Because the virtual-clock speed is constant between discrete changes, virtual time is a piecewise linear function of actual time, as depicted in Fig. 5(a), where t_s (speed change) is the latest speed change before arbitrary time t . The kernel keeps track of the most recent such actual time as `last_act`, the corresponding virtual time as `last_virt`, and the current speed of virtual time as `speed`. These values are initialized in `initialize()` and updated in `change_speed()`.

The convenience function `act_to_virt()` converts an actual time to a virtual time, assuming that `act > last_act` and that there is no virtual-clock speed change between `last_act` and `act`. By (4), the virtual-clock speed at t is the slope of the line graphed in Fig. 5(a) with $t_s = \text{last_act}$, resulting in the simple calculation performed in that function. Similarly, the convenience function `virt_to_act()` converts a virtual time to an actual time, assuming that `virt > last_virt` and that there is no virtual-clock speed change between `last_act` and `virt_to_act(virt)`.

In order to set the release timer for a level-C job, the kernel invokes `schedule_pending_release()`. This function uses `virt_to_act(v(r_{i,k}))` to determine when the timer should fire. This time could be incorrect if the

```

Function initialize()
1 | last_act := now();
2 | last_virt := 0;
3 | speed := 1;

Function act_to_virt (act)
4 | return last_virt + (act - last_act) · speed;

Function virt_to_act (virt)
5 | return last_act + (virt - last_virt)/speed;

Function schedule_pending_release ( $\tau_{i,k}$ ,
 $v(r_{i,k})$ )
6 | Set release timer to fire at virt_to_act ( $v(r_{i,k})$ );

Function job_release ( $\tau_{i,k}$ )
7 |  $r_{i,k} := \text{now}()$ ;
8 |  $v(y_{i,k}) := \text{act\_to\_virt}(r_{i,k}) + Y_i$ ;
9 |  $y_{i,k} := \perp$ ;

Function job_complete ( $\tau_{i,k}$ )
10 | virt := act_to_virt (now());
11 | if  $y_{i,k} = \perp$  and  $v(y_{i,k}) < \text{virt}$  then
12 |   |  $y_{i,k} := \text{virt\_to\_act}(v(y_{i,k}))$ ;
13 |   Report  $\tau_{i,k}$ ,  $r_{i,k}$ ,  $y_{i,k}$ , now(), and whether the
   | level-C ready queue is empty to the monitor program;

Function change_speed (new_speed)
14 | act := now();
15 | virt := act_to_virt (act);
16 | foreach  $\tau_{i,k}$  such that  $y_{i,k} = \perp$  and  $v(y_{i,k}) < \text{virt}$  do
17 |   |  $y_{i,k} := \text{virt\_to\_act}(v(y_{i,k}))$ ;
18 | last_act := act;
19 | last_virt := virt;
20 | speed := new_speed;
21 | foreach  $\tau_{i,k}$  such that a pending release has been
   | scheduled for virtual time  $v(r_{i,k})$  do
22 |   | Reset release timer to fire at
   |   | virt_to_act ( $v(r_{i,k})$ );

```

Algorithm 1: In-kernel functionality used to handle virtual time.

virtual-clock speed is changed before the timer fires, but in that case `change_speed()` will update the timer to fire at the correct time.

When a job release actually occurs, `job_release()` is called. This function determines the scheduling priority of $\tau_{i,k}$, which is simply the virtual time $v(y_{i,k})$ because the actual time $y_{i,k}$ is not known until $y_{i,k}$ occurs (because the virtual-clock speed may change). However, recall that the definition of “response-time tolerance” in Def. 1 is based on the actual time $y_{i,k}$. Therefore, it will generally be necessary for the kernel to determine $y_{i,k}$ and return it to the monitor program. Initially, the kernel uses the placeholder \perp , to indicate that $y_{i,k}$ has not yet occurred. There are three cases for when $y_{i,k}$ could occur relative to $t_{i,k}^c$, as depicted in Fig. 5(b)–(d).

If $t_{i,k}^c \leq y_{i,k}$, as depicted in Fig. 5(b), then $\tau_{i,k}$ meets its response-time tolerance (which was defined in Def. 1 to be nonnegative) by definition. Therefore, it is sufficient to

return \perp to the monitor program in this situation.

If $t_{i,k}^c > y_{i,k}$ and the speed of the virtual clock changes at least once between $y_{i,k}$ and $t_{i,k}^c$, then this scenario is depicted in Fig. 5(c), where t_s now refers to the first virtual-clock speed change *after* $y_{i,k}$. In this case, $y_{i,k}$ is computed when `change_speed()` is called at time t_s .

If $t_{i,k}^c > y_{i,k}$ and the speed of the virtual clock does not change between $y_{i,k}$ and $t_{i,k}^c$, then this scenario is depicted in Fig. 5(d). In this case, $y_{i,k}$ is computed when `job_complete($\tau_{i,k}$)` is called.

When any $\tau_{i,k}$ completes, `job_complete()` performs the just-mentioned check and notifies the monitor program.

When the monitor program requests a virtual-clock speed change, `change_speed()` is called. This function performs the updates mentioned above and also updates `last_act`, `last_virt`, and `speed` so that `virt_to_act()` and `act_to_virt()` remain correct.

The general structure of a userspace monitor program is presented in Algo. 2. A significant portion of the code is intended to detect the earliest possible idle normal instant. We define the following definition, which is closely related to the definition of “idle normal instant” in Def. 2.

Def. 3. t is a *candidate idle instant* at time $t_2 \geq t$ if some processor is idle at t and any job pending at t either meets its response-time tolerance or is still pending at t_2 .

In Fig. 4, t is a candidate idle instant at t_2 even if $\tau_{1,3}$ misses its response-time tolerance, as long as $\tau_{1,2}$ and $\tau_{2,5}$ meet their response-time tolerances.

The following theorem shows that we may consider only one candidate idle instant at any given time and still find the earliest idle normal instant. In Fig. 4, t_2 was selected as a time when a processor becomes idle, in order to illustrate this theorem.

Theorem 1. *If t is a candidate idle instant at t_2 and t_2 is an idle normal instant, then t is an idle normal instant.*

Proof. Because t is a candidate idle instant, by Def. 3, every job pending at t that is no longer pending at t_2 meets its response-time tolerance. Furthermore, because t_2 is an idle normal instant, by Def. 2, every job that is still pending at t_2 meets its response-time tolerance. Therefore, every job pending at t meets its response-time tolerance. Furthermore, because t is a candidate idle instant, by Def. 3, some processor is idle at t . Therefore, by Def. 2, t is an idle normal instant. \square

In order to detect an idle normal instant, we maintain as `idle_cand` the earliest candidate idle instant and as `pend_idle_cand` the set of incomplete jobs pending at `idle_cand`. If there is no current candidate idle instant, then the placeholder \perp is used instead. So that the monitor program can determine `pend_idle_cand` when a candidate idle instant is detected, it always maintains as `pend_now` the set of jobs currently pending. A job is added to `pend_now` whenever it is released, in `on_job_release()`. A job is removed from `pend_now` as soon as it has completed, in `on_job_complete()`.

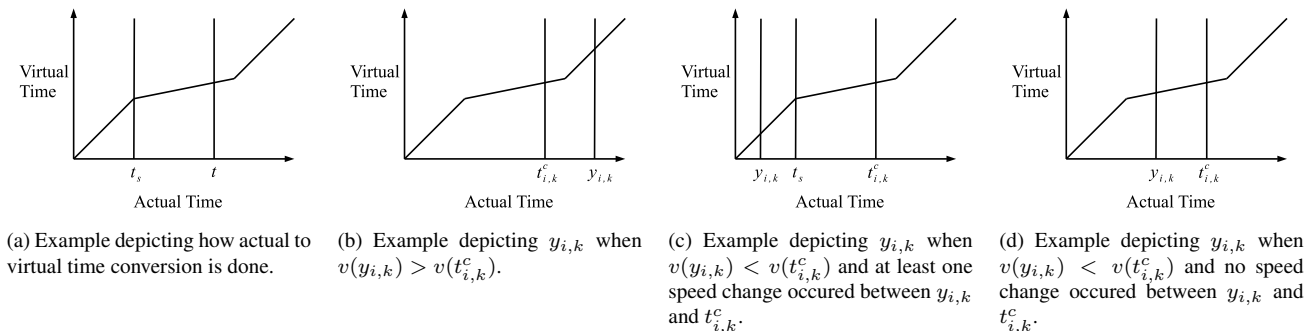


Figure 5: Examples illustrating virtual time computations in the kernel.

In the definition of `job_complete()` in Algo. 1, the kernel reports to the monitor program whether the ready queue is empty. The purpose for this reporting is that if the ready queue is empty, the processor that just completed $\tau_{i,k}$ has become idle. This fact is exploited in both `init_recovery()` and `on_job_complete()` in order to detect candidate idle instants.

The function `init_recovery()` initializes *recovery mode*, the process of finding an idle normal instant. Because recovery mode is always initiated as a result of a job missing its response-time tolerance, and such a miss is detected when the job completes, there is a relevant job completion time `comp_time`. As discussed above, if the ready queue was empty at `comp_time`, as indicated by `queue_empty`, then a processor became idle at `comp_time`. Therefore, by the definition of “candidate idle instant” in Def. 3, `comp_time` is a candidate idle instant. This case is handled in Lines 2–4. Otherwise, no candidate idle instant has yet been detected, as handled in Lines 5–7.

As discussed above, the function `on_job_release()` simply updates the set `pend_now` of currently pending jobs.

The function `on_job_complete()` first, in Line 9, updates `pend_now` as discussed above. Then, Lines 10–11 consider a response-time tolerance miss, regardless of whether the monitor program is currently in recovery mode. If such a miss occurs, then the function `handle_miss()` is called. The particular implementation of this function differs among the monitor programs we consider, and is discussed later. Lines 12–20 execute only when the monitor program is in recovery mode. Lines 12–17 execute if there is already a candidate idle instant under consideration. Lines 13–15 execute if $\tau_{i,k}$ has missed its response-time tolerance, in which case any prior candidate idle instant is no longer a candidate idle instant. On the other hand, Lines 16–17 execute when $\tau_{i,k}$ has met its response-time tolerance. In this case, we remove $\tau_{i,k}$ from the set `pend_idle_cand` of still-pending jobs that were pending at `idle_cand`. Lines 18–20 consider the case that `comp_time` has become the earliest candidate idle instant. This could happen either because a processor just became idle, or because a previous candidate idle instant was just discarded in Lines 13–15 while a processor was idle. In either case,

we start considering `comp_time` as a candidate idle instant. Finally, in Lines 21–23 we consider the case that there is an existing candidate idle instant, but the set `pend_idle_cand` is empty. Whenever this situation occurs, either because the last job in `pend_idle_cand` was removed on Line 17 or because the set of pending jobs considered in Line 20 was empty, the system exits recovery mode.

Our first userspace monitor program, `SIMPLE`, is depicted in Algo. 3. It is given the response-time tolerances desired for the tasks and a virtual time speed $s(t)$ used for overload recovery. When a response-time tolerance miss is detected while the system is not in recovery mode, it simply slows down the virtual clock and starts recovery mode.

Our second userspace monitor program, `ADAPTIVE`, is depicted in Algo. 4. It allows a value of $s(t)$ to be determined at runtime, selecting a smaller value for a more significant response-time tolerance miss. This minimizes the impact on the system when only a minor response-time tolerance miss has occurred, but provides a more drastic response when a larger miss has occurred. The monitor accepts an *aggressiveness factor* a in addition to the set of response-time tolerances, providing additional tuning. Once a response-time tolerance violation is detected, the monitor maintains the invariant that $s(t) = a \cdot \min((Y_i + \xi_i)/R_{i,k})$, where the `min` is over all jobs with $t_{i,k}^c$ after recovery mode last started. Thus, it chooses the speed based on the largest observed response time since recovery mode started.

5 Experiments

When a designer provisions an MC system, he or she should select level-C PWCETs that will be infrequently violated. Therefore, in the most common cases, overload conditions should be inherently transient, and it should be possible to return the system to normal operation relatively quickly. Therefore, our experiments consist of transient overloads rather than continuous overloads.

We ran experiments on a system with one quad-core 920-i7 CPU at 2.67 GHz, with 4GB of RAM. We generated 20 task sets, using a methodology similar to that described in [11], which used task systems designed to mimic avionics. We generated task systems where levels A and B each

```

Function init_recovery (comp_time,
queue_empty)
1  | recovery_mode := true;
2  | if queue_empty then
3  |   | idle_cand := comp_time;
4  |   | pend_idle_cand := pend_now;
5  | else
6  |   | idle_cand :=  $\perp$ ;
7  |   | pend_idle_cand := {};

Function on_job_release ( $\tau_{i,k}$ )
8  | Add  $\tau_{i,k}$  to pend_now;

Function on_job_complete ( $\tau_{i,k}, r_{i,k}, y_{i,k}$ ,
comp_time, queue_empty)
9  | Remove  $\tau_{i,k}$  from pend_now;
10 | if comp_time -  $y_{i,k} > \xi_i$  then
11 |   | handle_miss ( $\tau_{i,k}, r_{i,k}, y_{i,k}$ , comp_time,
12 |   | queue_empty);
13 | if recovery_mode and idle_cand  $\neq \perp$  then
14 |   | if comp_time -  $y_{i,k} > \xi_i$  then
15 |   |   | idle_cand :=  $\perp$ ;
16 |   |   | pend_idle_cand := {};
17 |   | else
18 |   |   | Remove  $\tau_{i,k}$  from pend_idle_cand;
19 |   | if recovery_mode and idle_cand =  $\perp$  and
20 |   | queue_empty then
21 |   |   | idle_cand := comp_time;
22 |   |   | pend_idle_cand := pend_now;
23 |   | if recovery_mode and idle_cand  $\neq \perp$  and
24 |   | pend_idle_cand = {} then
25 |   |   | change_speed (1);
26 |   |   | recovery_mode := false;

```

Algorithm 2: Userspace monitor algorithms common to SIMPLE and ADAPTIVE.

```

Function handle_miss ( $\tau_{i,k}, r_{i,k}, y_{i,k}$ , comp_time,
queue_empty)
1  | if not recovery_mode then
2  |   | change_speed ( $s(t)$ );
3  |   | init_recovery (comp_time,
4  |   | queue_empty);

```

Algorithm 3: Specific userspace implementation for SIMPLE.

occupy 5% of the system’s processor capacity and level C occupies 65% of the system’s capacity, *assuming that all jobs at all levels execute for their level-C PWCETs*. As in [11], we assumed that each task’s level-B PWCET is ten times its level-C PWCET, and that its level-A PWCET is twenty times its level-C PWCET.

At levels A and B, we generated tasks on one CPU at a time, using 5% of each CPU’s capacity for level A (assuming level-C execution times) and 5% for level B (again assuming level-C execution times). For level-A tasks, we selected periods randomly from the set {25 ms, 50 ms, 100 ms}. For level-B tasks, we selected ran-

```

Function handle_miss ( $\tau_{i,k}, r_{i,k}, y_{i,k}$ , comp_time,
queue_empty)
1  | if not recovery_mode then
2  |   | current_speed := 1;
3  |   | init_recovery (comp_time,
4  |   | queue_empty);
5  |   | new_speed :=  $a \cdot (Y_i + \xi_i) / (\text{comp\_time} - r_{i,k})$ ;
6  |   | if new_speed < current_speed then
7  |   |   | change_speed (new_speed);
8  |   |   | current_speed := new_speed;

```

Algorithm 4: Specific userspace implementation for ADAPTIVE.

dom multiples of the largest level-A period on the same CPU, capped at 300 ms. We then selected, for each task, a utilization (at its own criticality level) uniformly from (0.1, 0.4). This is the “uniform medium” distribution from prior work, e.g., [5]. For utilization at level C, the resulting choice is scaled by 1/20 for level-A tasks and 1/10 for level-B tasks. When a task would not fit within the allocated capacity for its criticality level, its utilization was scaled down to fit. Each task was then assigned a level-C PWCET based on multiplying its level-C utilization by its period.

At level C, we selected periods that were multiples of 5 ms between 10 ms and 100 ms, inclusive. We used uniform medium utilizations, as at levels A and B. As we did with levels A and B, we scaled down the utilization of the last task to fit. Y_i was selected for each level-C task using G-FL, which provides better response time bounds than G-EDF [9]. To determine response-time tolerances, we used the analytical bounds described in our technical report [8].

We tested the following overload scenarios:

- (SHORT) - All jobs at levels A, B, and C execute for their level-B PWCETs for 500 ms, and then execute for their level-C PWCETs afterward.
- (LONG) - All jobs at levels A, B, and C execute for their level-B PWCETs for 1 s, and then execute for their level-C PWCETs afterward.
- (DOUBLE) - All jobs at levels A, B, and C execute for their level-B PWCETs for 500 ms, execute for their level-C PWCETs for one second, execute for their level-B PWCETs for another 500 ms, and then execute for their level-C PWCETs afterward.

Because levels A and B together occupy 10% of the system’s capacity at level C, and because level-B PWCETs are ten times more pessimistic than level-C PWCETs, these represent a particularly pessimistic scenario in which all CPUs are occupied by level-A and -B work for almost all of the time during the overload.

For each overload scenario, we used SIMPLE with $s(t)$ choices from 0.2 to 1 in increments of 0.2. The choice of $s(t) = 1$ does not use our overload management techniques at all and provides a baseline for comparison. We also used ADAPTIVE with a choices from 0.2 to 1.0 in increments of

0.2. We then recorded the minimum virtual-time speed (to analyze the behavior of ADAPTIVE) and the amount of time from when the last overload stopped until the virtual-time clock was returned to normal. We then averaged each result over all generated task sets.

In Fig. 6, we depict the average dissipation time using SIMPLE with respect to the choice of $s(t)$ during recovery. Additionally, we depict error bars for 95% confidence intervals. Under LONG, dissipation times are approximately twice as long as under SHORT. This is to be expected, because overhead occurs for twice as long. Under DOUBLE, dissipation times are bigger than under SHORT for $s(t) = 1$, but nearly identical for smaller choices of $s(t)$. This occurs because dissipation time is measured from the end of the *second* (and final) interval during which overload occurs. For sufficiently small choices of $s(t)$, the system usually recovers completely before the second interval of overload starts, and that interval is the same length as in SHORT. In any case, a reduction of at least 50% of the dissipation time can be achieved with a choice of $s(t) = 0.6$, and with that choice, the dissipation time is less than twice the length of the interval during which overload occurs. Smaller choices of $s(t)$ have diminishing returns, with only a small improvement in dissipation time. Such a small improvement is likely outweighed by the larger impact on job releases from selecting a smaller $s(t)$.

In Fig. 7, we depict the average dissipation time using ADAPTIVE with respect to the aggressiveness factor. As before, we depict error bars for 95% confidence intervals. There is significant variance in the initial choice of $s(t)$ by ADAPTIVE, depending on which level-C jobs complete first after the overload starts, resulting in the larger confidence intervals. This effect is particularly pronounced in the case of DOUBLE. By comparing Figs. 6 and 7, we see that ADAPTIVE significantly reduces the dependency of dissipation time on the length of the overload interval. Furthermore, dissipation times are often significantly smaller under ADAPTIVE than under SIMPLE.

However, in order to fully evaluate ADAPTIVE, we must consider the minimum $s(t)$ value it chooses. Fig. 8 depicts the average of this choice with respect to the aggressiveness value, in addition to 95% confidence intervals. Here, we see that ADAPTIVE achieves smaller dissipation times than SIMPLE by choosing significantly slower virtual-clock speeds. Thus, jobs are released at a drastically lower frequency during the recovery period. Therefore, under the highly pessimistic scenarios we considered, SIMPLE is a better choice than ADAPTIVE.

As discussed above, level-C tasks run very little during the overload, so jobs pending at the end of the overload dominate other jobs in producing the largest response times. Because ADAPTIVE usually results in complete recovery from overload before the second overload interval, this causes nearly identical minimum choices of $s(t)$ between SHORT and DOUBLE. Similarly, because the overload interval is twice as long under LONG than under SHORT, the minimum choice of $s(t)$ is about half under LONG com-

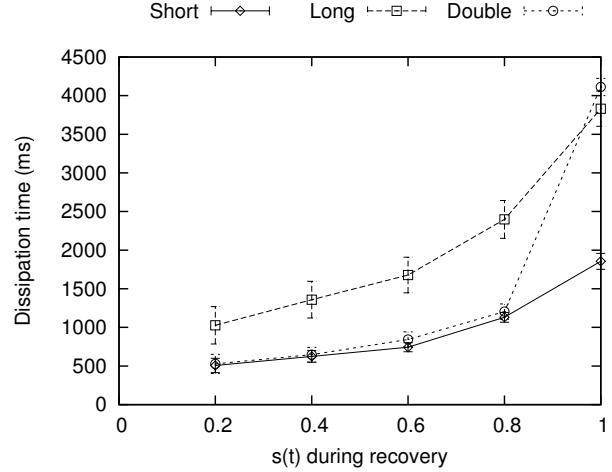


Figure 6: Dissipation time for SIMPLE

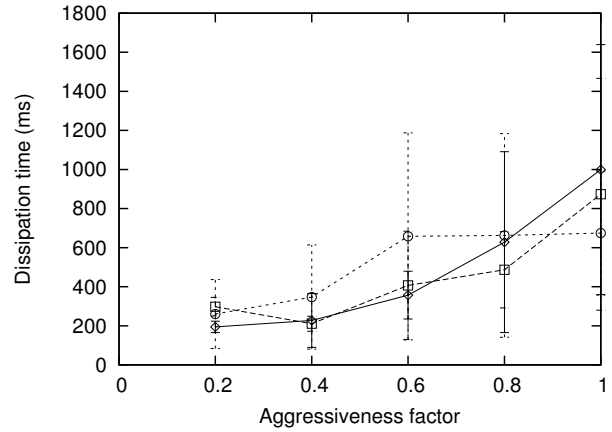


Figure 7: Dissipation time for ADAPTIVE

pared to SHORT.

In summary, the best choice of monitor under the tested conditions was SIMPLE with $s(t) = 0.6$, although $s(t) = 0.8$ could be a good choice if it is preferable to have a smaller impact on new releases with a longer dissipation time.

We also measured the same overheads considered in [11] both with and without our virtual time mechanism present, and considering both average and maximum observed overheads. For most overheads considered, there was no significant difference from the virtual time mechanism. However, there was variance in the scheduling overheads, as depicted in Fig. 9. For average-case overheads, the introduction of virtual time increased the scheduling time by about 40%, while for worst-case overheads, the introduction of virtual time approximately doubled the scheduling time. Because level C is SRT, the average-case overheads are more relevant, and the cost of adding the virtual time mechanism is small. Furthermore, the userspace monitor program had an effective CPU utilization of approximately 0.1, less than a typical task.

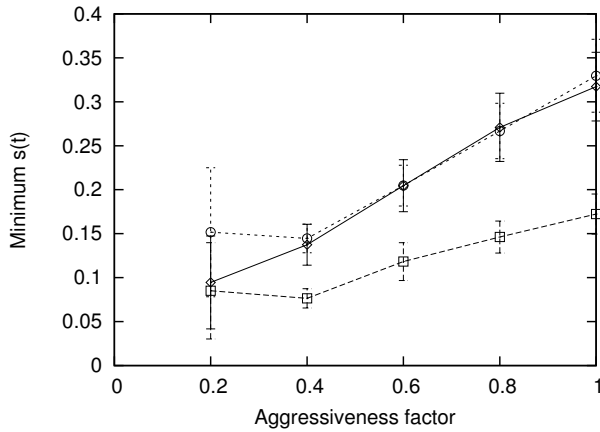


Figure 8: Minimum $s(t)$ for ADAPTIVE

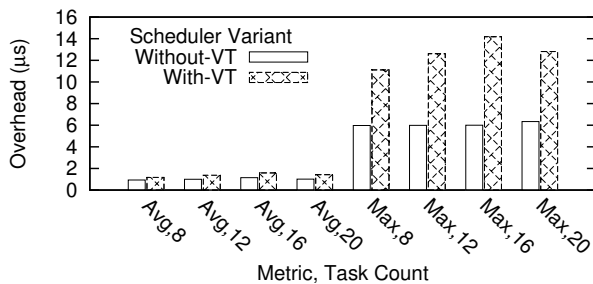


Figure 9: Scheduling overhead measurements

6 Conclusion

In this paper, we addressed the problem of scheduling under MC^2 when a transient overload occurs. We discussed the conditions that could cause an overload to result in a long-running increase in response-time bounds, and proposed a virtual-time mechanism to deal with these conditions.

We then presented an implementation of our mechanism and provided experiments to demonstrate that it can effectively provide recovery from unexpected overload scenarios. In our experiments, dissipation times could be brought within twice the length of a pessimistic overload scenario by only moderately affecting the time between job releases, and our scheme created little additional overhead.

References

- [1] LITMUS^{RT} home page. <http://www.litmus-rt.org/>.
- [2] S. Baruah, G. Koren, B. Mishra, A. Raghunathan, L. Rosier, and D. Shasha. On-line scheduling in the presence of overload. In *FOCS*, pages 100–110, 1991.
- [3] G. Beccari, S. Caselli, M. Reggiani, and F. Zanichelli. Rate modulation of soft real-time tasks in autonomous robot control systems. In *ECRTS*, pages 21–28, 1999.
- [4] A. Block. *Adaptive Multiprocessor Real-Time Systems*. PhD thesis, The University of North Carolina at Chapel Hill, 2008.
- [5] B. Brandenburg. *Scheduling and Locking in Multiprocessor*

- Real-Time Operating Systems*. PhD thesis, The University of North Carolina at Chapel Hill, 2011.
- [6] A. Burns and R. Davis. Mixed criticality systems - a review. <http://www-users.cs.york.ac.uk/~burns/review.pdf>, December 2013.
- [7] G. Buttazzo and J. Stankovic. RED: Robust earliest deadline scheduling. In *3rd International Workshop on Responsive Computing Systems*, pages 100–111, 1993.
- [8] J. Erickson and J. Anderson. Dissipation bounds: Recovering from overload in multicore mixed-criticality systems. Technical Report TR14-001, Department of Computer Science, University of North Carolina at Chapel Hill, Chapel Hill, NC, May 2014.
- [9] J. Erickson, J. Anderson, and B. Ward. Fair lateness scheduling: Reducing maximum lateness in G-EDF-like scheduling. *Real-Time Systems*, 50(1):5–47, 2014.
- [10] P. Gargali. On best-effort utility accrual real-time scheduling on multiprocessors. Master’s thesis, The Virginia Polytechnic Institute and State University, 2010.
- [11] J. Herman, C. Kenna, M. Mollison, J. Anderson, and D. Johnson. RTOS support for multicore mixed-criticality systems. In *RTAS*, pages 197–208, 2012.
- [12] M. Jan, L. Zaourar, and M. Pitel. Maximizing the execution rate of low criticality tasks in mixed criticality systems. In *WMC, RTSS*, pages 43–48, 2013.
- [13] G. Koren and D. Shasha. D^{over} : An optimal on-line scheduling algorithm for overloaded real-time systems. In *RTSS*, pages 290–299, 1992.
- [14] H. Leontyev and J. H. Anderson. Generalized tardiness bounds for global multiprocessor scheduling. *Real-Time Sys.*, 44(1):26–71, 2010.
- [15] H. Leontyev, S. Chakraborty, and J. Anderson. Multiprocessor extensions to real-time calculus. In *RTSS*, pages 410–421, 2009.
- [16] C. Locke. *Best-Effort Decision Making for Real-Time Scheduling*. PhD thesis, Carnegie Mellon University, 1986.
- [17] M. Mollison, J. Erickson, J. Anderson, S. Baruah, and J. Scoredos. Mixed-criticality real-time scheduling for multicore systems. In *ICSS*, pages 1864–1871, 2010.
- [18] F. Santy, L. George, P. Thierry, and J. Goossens. Relaxing mixed-criticality scheduling strictness for task sets scheduled with fp. In *ECRTS*, pages 155–165, 2012.
- [19] F. Santy, G. Raravi, G. Nelissen, V. Nélis, P. Kumar, J. Goossens, and E. Tovar. Two protocols to reduce the criticality level of multiprocessor mixed-criticality systems. In *RTNS*, pages 183–192, 2013.
- [20] I. Stoica, H. Abdel-Wahab, K. Jeffay, S. Baruah, J. Gehrke, and C. Plaxton. A proportional share resource allocation algorithm for real-time, time-shared systems. In *RTSS*, pages 288–299, 1996.
- [21] H. Su and D. Zhu. An elastic mixed-criticality task model and its scheduling algorithm. In *DATE*, pages 147–152, 2013.
- [22] H. Su, D. Zhu, and D. Mosse. Scheduling algorithms for elastic mixed-criticality tasks in multicore systems. In *RTCSA*, 2013.
- [23] L. Zhang. Virtual clock: A new traffic control algorithm for packet switching. In *SIGCOMM*, pages 19–29, 1990.