

# Resource-Sharing Servers for Open Environments

Marko Bertogna, Nathan Fisher, and Sanjoy Baruah

**Abstract**—We study the problem of executing a collection of independently designed and validated task systems upon a common platform composed of a preemptive processor and additional shared resources. We present an abstract formulation of the problem and identify the major issues that must be addressed in order to solve this problem. We present and prove the correctness of algorithms that address these issues, and thereby obtain a design for an open real-time environment.

**Index Terms**—Critical sections, earliest deadline first, hierarchical systems, open environments, resource-sharing systems, stack resource policy.

## I. INTRODUCTION

THE design and implementation of *open* real-time environments [20] is currently one of the more active research areas in the discipline of real-time computing. Such open environments aim to offer support for real-time *multiprogramming*: they permit multiple independently developed and validated real-time applications to execute concurrently upon a shared platform. That is, if an application is validated to meet its timing constraints when executing in isolation, then an open environment that accepts (or *admits*, through a process of admission control) this application guarantees that it will continue to meet its timing constraints upon the shared platform. The open environment has a runtime scheduler which arbitrates access to the platform among the various applications; each application has its own local scheduler for deciding which of its competing jobs executes each time the application is selected for execution by the “higher level” scheduler. (In recognition of this two-level scheduling hierarchy, such open environments are also often referred to as “hierarchical” real-time environments.)

In order to provide support for such real-time multiprogramming, open environments have typically found it necessary to place restrictions upon the structures of the individual applications. The first generation of such open platforms (see, e.g., [11], [17], [21], [22], [35], and [41])—this list is by no means exhaustive) assumed either that each application is composed of a finite collection of independent preemptive periodic (Liu and Layland) tasks [31], or that each application’s schedule is statically

precomputed and runtime scheduling is done via table lookup. Furthermore, these open environments focused primarily upon the scheduling of a single (fully preemptive) processor, ignoring the fact that runtime platforms typically include additional resources that may not be fully preemptable. The few [16], [19], [24], [38] that do allow for such additional shared resources typically make further simplifying assumptions on the task model, e.g., by assuming that the computational demands of each application may be aggregated and represented as a single periodic task, excluding the possibility to address hierarchical systems.

More recently, researchers have begun working upon open environments that are capable of operating upon more complex platforms. There are recent publications proposing designs for open environments that allow for sharing other resources in addition to the preemptive processor [9], [10], [18], [40]. These designs assume that each individual application may be characterized as a collection of sporadic tasks [8], [32], distinguishing between shared resources that are *local* to an application (i.e., only shared within the application) and *global* (i.e., may be shared among different applications). However, these approaches either propose that global resources be executed with local preemptions disabled [18], potentially causing intolerable blocking inside an application; or allow applications to overrun their budget, while holding a lock [9], [10], [18], [40], increasing in this way the blocking among applications.

In this paper, we describe our design of an open environment upon a computing platform composed of a single preemptive processor and additional shared resources. We assume that each application can be modeled as a collection of preemptive jobs which may access shared resources within critical sections. (Such jobs may be generated by, for example, periodic and sporadic tasks.) We require that each such application be scheduled using some local scheduling algorithm, with resource contention arbitrated using some strategy such as the Stack Resource Policy (SRP) [3]. We describe what kinds of analysis such applications must be subject to and what properties these applications must satisfy, in order for us to be able to guarantee that they will meet their deadlines in the open environment.

The remainder of this paper is organized as follows. The rationale and design of our open environment is described in Sections II and III. In Section II, we provide a high-level overview of our design, and detail the manner in which we expect individual applications to be characterized—this characterization represents the *interface* specification between the open environment and individual applications running on it—and in Section III, we present the scheduling and admission-control algorithms used by our open environment. In Section IV, we discuss what kind of applications may be scheduled by our open environment. Section V analyzes the local feasibility problem of an admitted application. Section VI provides further details on one important interface parameter

Manuscript received December 18, 2008; revised May 31, 2009. First published July 21, 2009; current version published August 07, 2009. This work was supported in part by the NSF under Grant CNS-0834270, Grant CNS-0834132, Grant CNS-0615197, and ARO Grant W911NF-06-1-0425 and in part by the Wayne State University Faculty Research Award. Paper no. TII-08-12-0226.R1.

M. Bertogna is with Scuola Superiore Sant’Anna, Pisa 56124, Italy (e-mail: marko@sssup.it).

N. Fisher is with the Department of Computer Science, Wayne State University, Detroit, MI 48202 USA (e-mail: fishern@cs.wayne.edu).

S. Baruah is with the Department of Computer Science, University of North Carolina at Chapel Hill, Chapel Hill, NC 27599-3175 USA (e-mail: baruah@cs.unc.edu).

Digital Object Identifier 10.1109/TII.2009.2026051

U.S. Government work not protected by U.S. copyright.

related to the locking of shared resources. In Section VII, we discuss possible design choices that can improve the local and global schedulability of applications scheduled in our open environment. In Section VIII, we relate our open environment framework to other previously proposed works. Finally, in Section IX, we present our conclusions.

## II. SYSTEM MODEL

In an open environment, there is a shared processing platform upon which several independent applications  $A_1, \dots, A_q$  execute. We also assume that the shared processing platform is composed of a single preemptive processor (without loss of generality, we will assume that this processor has unit computing capacity), and  $m$  additional (global) shared resources which may be shared among the different applications. Each application may have additional “local” shared logical resources that are shared between different jobs within the application itself—the presence of these local shared resources is not relevant to the design and analysis of the open environment. We will distinguish between:

- a unique *system-level scheduler* (or *global scheduler*), which is responsible for scheduling all admitted applications on the shared processor;
- one or more *application-level schedulers* (or *local schedulers*), that decide how to schedule the jobs of an application.

An *interface* must be specified between each application and the open environment. The goal of this interface specification is to abstract out and encapsulate the salient features of the application’s resource requirements. The open environment uses this information during *admission control*, to determine whether the application can be supported concurrently with other already admitted applications; for admitted applications, this information is also used by the open environment during runtime to make scheduling decisions. If an application is admitted, the interface represents its “contract” with the open environment, which may use this information to enforce (“police”) the application’s runtime behavior. As long as the application behaves as specified by its interface, it is guaranteed to meet its timing constraints; if it violates its interface, it may be penalized while other applications are isolated from the effects of this misbehavior. We require that the interface for each application  $A_k$  be characterized by three parameters.

- A *virtual processor* (VP) *speed*  $\alpha_k$ .
- A *jitter tolerance*  $\Delta_k$ .
- For each global shared resource  $R_\ell$ , a *resource holding time*  $H_k(R_\ell)$ .

The intended interpretation of these interface parameters is as follows: *all jobs of the application will meet their deadlines if executing upon a processor of computing capacity  $\alpha_k$ , with a service delay of at most  $\Delta_k$  time-units, and will lock resource  $R_\ell$  for no more than  $H_k$  time-units at a time during such execution.*

We now provide a brief overview of the application interface parameters. Section VI provides a more in depth discussion of the resource holding time parameter.

1) *VP Speed  $\alpha_k$* : Since each application  $A_k$  is assumed validated (i.e., analyzed for schedulability) upon a slower VP, this

parameter is essentially the computing capacity of the slower processor upon which the application was validated.

2) *Jitter Tolerance  $\Delta_k$* : Given a processor with computing capacity  $\alpha_k$  upon which an application  $A_k$  is validated, this is the maximum service delay that  $A_k$  can withstand without missing any of its deadlines. In other words,  $\Delta_k$  is the maximum release delay that all jobs can experience without missing any deadline.

At first glance, this characterization may seem like a severe restriction, in the sense that one will be required to “waste” a significant fraction of the VP’s computing capacity in order to meet this requirement. However, this is not necessarily correct. Consider the following simple (contrived) example. Let us represent a sporadic task [8], [32] by a three-tuple: (*WCET, relative deadline, period*). Consider the example application composed of the two sporadic tasks  $\{(1, 4, 4), (1, 6, 4)\}$  to be validated upon a dedicated processor of computing capacity one-half. The task set fully utilizes the VP. However, we could schedule this application such that no deadline is missed even when all jobs are released with a delay of two time units. That is, this application can be characterized by the pair of parameters  $\alpha_k = 1/2$  and  $\Delta_k = 2$ .

Observe that there is a correlation between the VP speed parameter  $\alpha_k$  and the timeliness constraint  $\Delta_k$ —increasing  $\alpha_k$  (executing an application on a faster VP) may cause an increase in the value of  $\Delta_k$ . Equivalently, a lower  $\alpha_k$  may result in a tighter jitter tolerance, with some job finishing close to its deadline. However, this relationship between  $\alpha_k$  and  $\Delta_k$  is not linear nor straightforward—by careful analysis of specific systems, a significant increase in  $\Delta_k$  may sometimes be obtained for a relatively small increase in  $\alpha_k$ .

Note that the validation process to derive the interface parameters for an application  $A_k$  does not require to effectively execute the application on a real processor of smaller speed. Such parameters might be derived from the worst-case execution times of the composing jobs on the (unit-speed) target platform, adopting proper schedulability tests associated to the local scheduling algorithm in use. More details on the validation process will be provided in Section IV.

Our characterization of an application’s processor demands by the parameters  $\alpha_k$  and  $\Delta_k$  is identical to the *bounded-delay resource partition* characterization of Feng and Mok [21], [22], [35] with the exception of the  $H_k(R_\ell)$  parameter.

3) *Resource Holding Times  $H_k(R_\ell)$* : For open environments which choose to execute all global resources disabling local preemptions (such as the design proposed in [18]),  $H_k(R_\ell)$  is simply the worst-case execution time upon the VP of the longest critical section holding global resource  $R_\ell$ . We have recently [12], [13], [23] derived algorithms for computing resource holding times when more general resource-access strategies such as the SRP [3] and the Priority Ceiling Protocol (PCP) [37], [39] are instead used to arbitrate access to these global resources; in [12], [13], and [23], we also discuss the issue of designing the specific application systems such that the resource holding times are decreased without compromising feasibility. We believe that our consideration of global shared resources—their abstraction by the  $H_k$  parameters in the interface, and the use we make of this information—is one of

Symbol	Description
$A_k$	$k$ -th application
$\alpha_k$	VP speed of $A_k$
$\Delta_k$	Jitter tolerance of $A_k$
$R_\ell$	$\ell$ -th global resource
$H_k(R_\ell)$	$A_k$ 's resource holding time for $R_\ell$
$\Pi(R_\ell)$	Global ceiling of resource $R_\ell$
$P_k$	Period of the server associated to $A_k$
$D_k$	Deadline of the server
$Z_k$	Reactivation time of the server
$V_k$	Virtual time of the server
$t_{cur}$	Current time instant
$J_{k,i}$	$i$ -th chunk of application $A_k$
$r(J_{k,i})$	Release time of $J_{k,i}$
$g(J_{k,i})$	Termination time of $J_{k,i}$
$d(J_{k,i})$	Deadline of $J_{k,i}$
$n_k$	Number of tasks composing the application $A_k$
$\tau_i$	$i$ -th task of the application
$c_i$	WCET of $\tau_i$
$d_i$	Deadline of $\tau_i$
$p_i$	Period of minimum inter-arrival time of $\tau_i$
$h_i(R_\ell)$	$\tau_i$ 's largest critical section on $R_\ell$
$\pi(R_\ell)$	Local ceiling of resource $R_\ell$

Fig. 1. Notation used throughout this paper.

our major contributions, and serves to distinguish our work from other projects addressing similar topics. Our approach toward resource holding times is discussed in greater detail in Section VI.

The notation adopted in this paper is summarized in Fig. 1.

### III. ALGORITHMS

In this section, we present the algorithms used by our open environment to make admission-control and scheduling decisions. We assume that each application is characterized by the interface parameters described in Section II. When a new application wishes to execute, it presents its interface to the *admission control algorithm*, which determines, based upon the interface parameters of this and previously admitted applications, whether to admit this application or not. If admitted, each application is executed through a dedicated server. At each instant during runtime, the (system-level) *scheduling algorithm* decides which server (i.e., application) gets to run. If an application violates the contract implicit in its interface, an *enforcement algorithm* polices the application—such policing may affect the performance of the misbehaving application, but should not compromise the behavior of other applications.

We first describe the global scheduling algorithm used by our open environment, in Section III-A. A description and proof of correctness of our admission control algorithm follows (in Sections III-B–III-E). The local schedulers that may be used by the individual applications will be addressed in Section IV.

#### A. System-Level Scheduler

Our scheduling algorithm is essentially an application of the Constant Bandwidth Server (CBS) of Abeni and Buttazzo [1], enhanced to allow for the sharing of nonpreemptable serially reusable resources and for the concurrent execution of different applications in an open environment. In the remainder of this paper, we will refer to this server with the acronym BROE: Bounded-Delay Resource Open Environment.

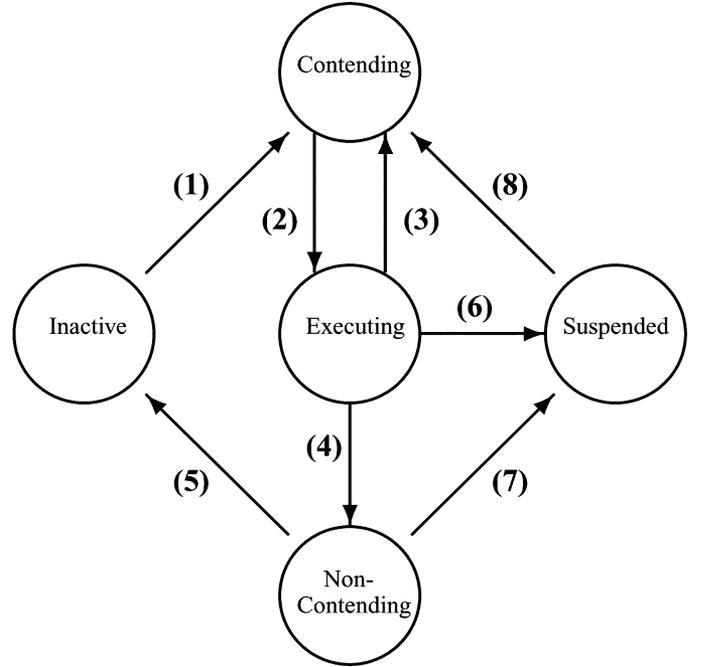


Fig. 2. State transition diagram. The labels on the nodes and edges denote the name by which the respective states and transitions are referred to in this paper.

CBS-like servers have an associated *period*  $P_k$ , reflecting the time-interval at which a continuously active server replenishes its budget. For a BROE server, the value assigned to  $P_k$  is as follows:

$$P_k \leftarrow \frac{\Delta_k}{2(1 - \alpha_k)}. \quad (1)$$

In addition, each server maintains three variables: a *deadline*  $D_k$ , a *virtual time*  $V_k$ , and a *reactivation time*  $Z_k$ . Since each application has a dedicated server, we will not make any distinction between server and application parameters. At each instant during runtime, each server assigns a *state* to the admitted application. There are five possible states (see Fig. 2). Let us define an application to be *backlogged* at a given time-instant if it has any active jobs awaiting execution at that instant, and *nonbacklogged* otherwise.

- Each nonbacklogged application is in either the *Inactive* or *Noncontending* states. If an application has executed for more than its “fair share,” then it is *Noncontending*; else, it is *Inactive*.
- Each backlogged application is in one the *Contending*, *Executing*, or *Suspended* states.<sup>1</sup> While contending, it is eligible to execute; executing for more than it is eligible results in it being suspended.

These variables are updated by BROE according to the following rules (i)–(ix) (let  $t_{cur}$  denote the current time).

- Initially, each application is in the *Inactive* state. If application  $A_k$  wishes to contend for execution at time-instant  $t_{cur}$  then it transits to the *Contending* state [transition

<sup>1</sup>Note that there is no analog of the *Suspended* state in the original definition of CBS [1].

(1) in Fig. 2]. This transition is accompanied by the following actions:

$$\begin{aligned} D_k &\leftarrow t_{\text{cur}} + P_k \\ V_k &\leftarrow t_{\text{cur}}. \end{aligned}$$

- (ii) At each instant, the system-level scheduling algorithm selects for execution some application  $A_k$  in the *Contending* state—the specific manner in which this selection is made is discussed below. Hence, observe that *only applications in the Contending state are eligible to execute*. An application scheduled for execution undergoes transition (2). The virtual time of an executing application  $A_k$  is incremented by the corresponding server at a rate  $1/\alpha_k$

$$\frac{d}{dt}V_k = \begin{cases} 1/\alpha_k, & \text{while } A_k \text{ is executing} \\ 0, & \text{the rest of the time} \end{cases}.$$

- (iii) When an application is preempted by a higher priority one, it undergoes transition (3) to the *Contending* state.  
 (iv) An application  $A_k$  which no longer desires to contend for execution (i.e., the application is no longer backlogged) transits to the *Noncontending* state [transition (4)], and remains there as long as  $V_k$  exceeds the current time.  
 (v) When  $t_{\text{cur}} \geq V_k$  for some such application  $A_k$  in the *Noncontending* state,  $A_k$  transits back to the *Inactive* state [transition (5)].  
 (vi) If the virtual time  $V_k$  of the executing application  $A_k$  becomes equal to  $D_k$ , then application  $A_k$  undergoes transition (6) to the *Suspended* state. This transition is accompanied by the following actions:

$$Z_k \leftarrow V_k \quad (2)$$

$$D_k \leftarrow V_k + P_k \quad (3)$$

- (vii) An application  $A_k$  in *Noncontending* state which desires to once again contend for execution (note  $t_{\text{cur}} < V_k$ , otherwise it would be in the *Inactive* state) transits to the *Suspended* state [transition (7)]. This transition is accompanied by the same actions [(2) and (3)] of transition (6).  
 (viii) An application  $A_k$  that is in the *Suspended* state necessarily satisfies  $Z_k \geq t_{\text{cur}}$ . As the current time  $t_{\text{cur}}$  increases, it eventually becomes the case that  $Z_k = t_{\text{cur}}$ . At that instant, application  $A_k$  transits back to the *Contending* state [transition (8)].  
 (ix) An application that wishes to gain access to a shared global resource  $R_\ell$  must perform a *budget check* (i.e., is there enough execution budget to complete execution of the resource prior to  $D_k$ ?). If  $\alpha_k(D_k - V_k) \geq H_k(R_\ell)$ , there is sufficient budget, and the server is granted access to resource  $R_\ell$ . Otherwise, the budget is insufficient to complete access to resource  $R_\ell$  by  $D_k$ . In this case, transition (6) is undertaken by an executing application immediately prior to entering an outermost critical section locking a global resource<sup>2</sup>  $R_\ell$ , updating the server parameters according to (2) and (3).

<sup>2</sup>Each application may have additional resources that are local in the sense that are not shared outside the application. Attempting to lock such a resource does not trigger transition (6).

Rules (i)–(viii) basically describe a bounded-delay version of the Constant Bandwidth Server, i.e., a CBS in which the maximum service delay experienced by an application  $A_k$  is bounded by  $\Delta_k$ . A similar server has also been used in [27] and [28]. The only difference from a straightforward implementation of a bounded-delay CBS is the deadline update of rule (vii) associated to transition (7) (which has been introduced in order to guarantee that when an application resumes execution, its relative deadline is equal to the server period, so that the budget is full) and the addition of rule (ix).

Rule (ix) has been added to deal with the problem of *budget exhaustion* when a shared resource is locked. This problem, previously described in [16] and [18], arises when an application accesses a shared resource and runs out of budget [i.e., is suspended after taking transition (2)] before being able to unlock the resource. This would cause intolerable blocking to other applications waiting for the same lock. If there is insufficient current budget, taking transition (6) right before an application  $A_k$  locks a critical section ensures that when  $A_k$  goes to the *Contending* state [through transition (8)], it will have  $D_k - V_k = P_k$ . This guarantees that  $A_k$  will receive  $(\alpha_k P_k)$  units of execution prior to needing to be suspended [through transition (6)]. Thus, *ensuring that the WCET of each critical section of  $A_k$  is no more than  $\alpha_k P_k$  is sufficient to guarantee that  $A_k$  experiences no deadline-postponement within any critical section*. Our admission control algorithm (Section III-B) does in fact ensure that

$$H_k(R_\ell) \leq \alpha_k P_k \quad (4)$$

for all applications  $A_k$  and all resources  $R_\ell$ ; hence, no lock-holding application experiences deadline postponement.

At first glance, requiring that applications satisfy Condition (4) may seem to be a severe limitation of our framework. But this restriction appears to be unavoidable if CBS-like approaches are used as the system-level scheduler: in essence, this restriction arises from a requirement that an application not get suspended (due to having exhausted its current execution capacity) whilst holding a resource lock. To our knowledge, all lock-based multilevel scheduling frameworks impose this restriction explicitly (e.g., [16]) or implicitly, by allowing lock-holding applications to continue executing nonpreemptively even when their current execution capacities are exhausted (e.g., [10], [18], and [40]).

We now describe how our scheduling algorithm determines which BROE server (i.e., which of the applications currently in the *Contending* state) to select for execution at each instant in time.

In brief, we implement EDF among the various contending applications, with the application deadlines (the  $D_k$ 's) being the deadlines under comparison. Access to the global shared resources is arbitrated using SRP.<sup>3</sup>

In greater detail:

- 1) Each global resource  $R_\ell$  is assigned a ceiling  $\Pi(R_\ell)$  which is equal to the minimum value from among all the period parameters  $P_k$  of  $A_k$  that use this resource.

<sup>3</sup>Recall that in our scheduling scheme, *deadline postponement cannot occur for an application while it is in a critical section*—this property is essential to our being able to apply SRP for arbitrating access to shared resources.

Initially,  $\Pi(R_\ell) \leftarrow \infty$  for all the resources. When an application  $A_k$  is admitted that uses global resource  $R_\ell$ ,  $\Pi(R_\ell) \leftarrow \min(\Pi(R_\ell), P_k)$ ;  $\Pi(R_\ell)$  must subsequently be recomputed when such an application leaves the environment.

- 2) At each instant, there is a *system ceiling*, which is equal to the minimum ceiling of any resource that is locked at that instant.
- 3) At the instant that an application  $A_k$  becomes the earliest-deadline one that is in the *Contending* state, it is selected for execution if and only if its period parameter  $P_k$  is strictly less than the system ceiling at that instant. Else, it is blocked while the currently executing application continues to execute.

As stated above, this is essentially an implementation of EDF+SRP among the applications. The SRP requires that the relative deadline of a job locking a resource be known beforehand; that is why our algorithm requires that deadline postponement not occur while an application has locked a resource.

### B. Admission Control

The admission control algorithm checks for four things.

- 1) As stated in Section III-A, we require that each application  $A_k$  have all its resource holding times (the  $H_k(R_\ell)$ 's) be  $\leq \alpha_k P_k$ —any application  $A_k$  whose interface does not satisfy this condition is summarily rejected. If the application is rejected, the designer may attempt to increase the  $\alpha_k$  parameter and resubmit the application; increasing  $\alpha_k$  will simultaneously increase  $\alpha_k P_k$ , while decreasing the  $H_k(R_\ell)$ 's.
- 2) The sum of the VP speeds—the  $\alpha_i$  parameters—of all admitted tasks may not exceed the computing capacity of the shared processor (assumed to be equal to one). Hence,  $A_k$  is rejected if admitting it would cause the sum of the  $\alpha_i$  parameters of all admitted applications to exceed one.
- 3) The effect of *interapplication blocking* must be considered—can such blocking cause any server to miss a deadline? A server-deadline miss occurs for a backlogged server when  $t_{\text{cur}} \geq D_k$  and  $V_k < D_k$ . The issue of interapplication blocking is discussed in the remainder of this section.
- 4) The admission of an application  $A_k$  (if it satisfies the above three items) must be delayed until the first time instant during which all global resources needed by  $A_k$  are available (i.e., not locked by other applications). This delay is necessary to avoid raising a resource's priority ceiling while it is currently locked by an application with larger period than  $A_k$ 's. As discussed in Theorem 1, the delay will help avoid undesirable blocking within SRP. Please note that the delay of application  $A_k$ 's admittance depends upon the system reaching a lock “stasis.” The length of this delay is bounded by the longest response time (elapsed total time between resource lock and release) of any critical section that accesses a resource needed by  $A_k$ . Thus,  $A_k$  is potentially delayed for a longer than usual amount of time; however, this decision prevents previously admitted applications from being penalized with additional blocking and avoids deviation from the standard SRP policy.

Admission control and *feasibility*—the ability to meet all deadlines—are two sides of the same coin. As stated above, our system-level scheduling algorithm is essentially EDF, with access to shared resources arbitrated by the SRP. Hence, the admission control algorithm needs to ensure that all the admitted applications together are feasible under EDF+SRP scheduling. We therefore looked to the EDF+SRP feasibility test in [6], [29], [36] for inspiration and ideas. In designing an admission control algorithm based upon these known EDF+SRP feasibility tests there are a series of design decisions. Based upon the available choices, we came up with two possible admission control algorithms: a more accurate one that requires information regarding each application's resource holding time for every resource, and a slightly less accurate test that reduces the amount of information required by the system to make an admission control decision.

Prior to introducing the admission control algorithms, Section III-C will prove that many of the desirable properties of SRP that hold for sporadic task systems [3] continue to hold for our BROE server. Section III-D will provide useful bounds on the demand of a server. Finally, Section III-E will describe and prove the correctness of two admission control algorithms.

### C. Stack-Resource Policy Properties

As mentioned at the beginning of this section,  $H_k(R_\ell) \leq \alpha_k P_k$  for every global resource used by application  $A_k$ . The previous considerations allow deriving some important properties for the open environment, since there won't be any deadline postponement inside a critical section, we can view each application execution as a release sequence of “chunks” (i.e., separate jobs), as suggested in [16]. A new chunk is released each time the application enters the *Contending* state [transitions (1) and (8)] and is terminated as soon as the state transitions from *Executing* [transitions (4) and (6)]. We will denote the  $\ell$ 'th chunk of application  $A_k$  as  $J_{k,\ell}$ . The *release time* of  $J_{k,\ell}$  is denoted as  $r(J_{k,\ell})$ . The *termination time* of  $J_{k,\ell}$  is denoted by  $g(J_{k,\ell})$ . Finally, the *deadline* of chunk  $J_{k,\ell}$  is the  $D_k$  value of the server at the time it transitioned to contending; the deadline of chunk  $J_{k,\ell}$  is represented by  $d(J_{k,\ell})$ . Let  $V_k(t)$  denote the server's value of  $V_k$  at time  $t$ .

A *priority inversion* between applications is said to occur during runtime if the earliest-deadline application that is contending—awaiting execution—at that time cannot execute because some resource needed for its execution is held by some other application. This (later-deadline) application is said to *block* the earliest-deadline application. SRP bounds the amount of time that any application chunk may be blocked. The enforcement mechanism used in our open environment allows proving the following Theorem.

*Theorem 1 (SRP Properties):* There are no deadlocks between applications in the open environment. Moreover, all chunks  $J_{k,\ell}$  of an application  $A_k$  that does not exceed the declared resource-holding-time have the following properties.

- $J_{k,\ell}$  cannot be blocked after it begins execution.
- $J_{k,\ell}$  may be blocked by at most one later deadline application for at most the duration of one resource-holding-time.

*Proof:* Note that admission-control property 4 of Section III-B ensures that resource ceilings are not raised

while an application is locked. Thus, our resource-arbitration protocol is identical to SRP and the proof is identical to the Proof of [3, Th. 6]. The only difference is that in our case the items to be scheduled are application chunks instead of jobs. ■

#### D. Bounding the Demand of Server Chunks

It is useful to quantify the amount of execution that a chunk of a server requires over any given time interval. We quantify the *demand* of a server chunk, and attempt to bound the total demand (over an interval of time) by a server for  $A_k$ . The bound on demand will be useful in the next subsection which discusses our admission control algorithms. The following are formal definitions of demand for a server chunk and the total demand for a server.

*Definition 1 (Demand of Server Chunk  $J_{k,\ell}$ ):* The *demand*  $W(J_{k,\ell}, t_1, t_2)$  of server chunk  $J_{k,\ell}$  over the interval  $[t_1, t_2]$  is the amount of execution that  $J_{k,\ell}$  (with deadline and release time in the interval  $[t_1, t_2]$ ) must receive before making a transition from *Executing* to *Noncontending* or *Suspended*. Formally,  $W(J_{k,\ell}, t_1, t_2)$  is equal to

$$\begin{cases} \alpha_k(V_k(g(J_{k,\ell})) - V_k(r(J_{k,\ell}))), & \text{if } (r(J_{k,\ell}) \geq t_1) \wedge \\ & (g(J_{k,\ell}) < d(J_{k,\ell}) \leq t_2); \\ \alpha_k P_k, & \text{if } (r(J_{k,\ell}) \geq t_1) \\ & \wedge (d(J_{k,\ell}) \leq t_2) \\ & \wedge (g(J_{k,\ell}) \geq d(J_{k,\ell})); \\ 0, & \text{otherwise.} \end{cases} \quad (5)$$

*Definition 2 (Cumulative Demand of BROE Server for  $A_k$ ):* The *cumulative demand* of  $A_k$  over the interval  $[t_1, t_2]$  is the total demand of all of  $A_k$ 's server chunks with both release times and deadlines within the interval  $[t_1, t_2]$

$$W(A_k, t_1, t_2) \stackrel{\text{def}}{=} \sum_{\ell \geq 1} W(J_{k,\ell}, t_1, t_2). \quad (6)$$

Different chunks of the same server may execute for different amounts of time. The reason is that some chunks may terminate early due to becoming noncontending or trying to enter a critical section [i.e., rule (iv) or (ix)]. For these chunks, the execution they receive may be less than  $\alpha_k P_k$ . Unfortunately, there are infinitely many possible application execution scenarios over any given interval (resulting in different sequences of state transitions). With all these possibilities, how does one determine the cumulative demand of  $A_k$  over any interval? Fortunately, we may, in fact, derive upper bounds for the cumulative demand of a server for specific sequences of chunks. The upper bound for these sequences will be used in proof of correctness for the admission control algorithm. In the remainder of this subsection, we will present a series of lemmas that derives the upper bound on the cumulative demand of a sequence of server chunks.

The first lemma states that the virtual time  $V_k$  of a server cannot exceed the deadline parameter  $D_k$ .

*Lemma 1:* For all chunks  $J_{k,\ell}$  of BROE server of  $A_k$

$$V_k(g(J_{k,\ell})) \leq d(J_{k,\ell}). \quad (7)$$

*Proof:* Observe that rule (vi) implies that the virtual time  $V_k$  does not exceed the current server deadline  $D_k$ . Therefore,

whenever any chunk  $J_{k,\ell}$  is terminated [via transitions (4), or (6)] at time  $g(J_{k,\ell})$ , the server's virtual time  $V_k(g(J_{k,\ell}))$  does not exceed the deadline  $d(J_{k,\ell})$  of the chunk. ■

The next lemma formally states that, when a chunk terminates with transition (6), the virtual time does not increase until the release of the next chunk.

*Lemma 2:* If  $J_{k,\ell+1}$  was released due to transition (8) (i.e., *Suspended* to *Contending*), then  $V_k(r(J_{k,\ell+1})) = V_k(g(J_{k,\ell}))$ .

*Proof:* If  $J_{k,\ell+1}$  was released due to transition (8), the transition prior to (8) must have been either (6), or the successive transitions of (4) and (7). For these transitions, one of the server rules (iv), (vi), (vii), or (ix) applies when terminating the previous chunk  $J_{k,\ell}$ . However, notice that none of these rules updates  $V_k$ , and since virtual time cannot progress unless the server is contending, the virtual time at the release of  $J_{k,\ell+1}$  [i.e.,  $V_k(r(J_{k,\ell+1}))$ ] must equal the virtual time at the termination of  $J_{k,\ell}$  [i.e.,  $V_k(g(J_{k,\ell}))$ ]. ■

In the final lemma of this subsection, we consider any sequence of chunks where the server does not become inactive in between releases and the virtual time at the release of the first chunk equals actual time. For such a sequence of chunks, we show that the demand of the chunks from the release time of the first chunk of the sequence to the deadline of the last chunk of the sequence does not exceed  $\alpha_k$  times the sequence length (i.e., the deadline of the last chunk minus release time of the first chunk).

*Lemma 3:* If  $J_{k,\ell}, J_{k,\ell+1}, \dots, J_{k,s}$  is a sequence of successively released chunks by the BROE server for  $A_k$  where  $J_{k,\ell}$  satisfies  $V_k(r(J_{k,\ell})) = r(J_{k,\ell})$ , and  $J_{k,\ell+1}, \dots, J_{k,s}$  were all released due to transition (8); if  $J_{k,\ell}, \dots, J_{k,s-1}$  meet their deadline, then

$$W(A_k, r(J_{k,\ell}), d(J_{k,s})) \leq \alpha_k(d(J_{k,s}) - r(J_{k,\ell})). \quad (8)$$

*Proof:* According to Definition 2,  $W(A_k, r(J_{k,\ell}), d(J_{k,s}))$  is the sum of the  $W(J_{k,i}, r(J_{k,\ell}), d(J_{k,s}))$  for each chunk  $J_{k,i}$ , where  $\ell \leq i \leq s$ . Since both  $r(J_{k,i})$  and  $d(J_{k,i})$  must be included in the interval  $[r(J_{k,\ell}), d(J_{k,s})]$  and chunks  $J_{k,\ell}, \dots, J_{k,s-1}$  meet their deadlines, (5) implies

$$\begin{aligned} W(A_k, r(J_{k,\ell}), d(J_{k,s})) &= W(J_{k,s}, r(J_{k,\ell}), d(J_{k,s})) \\ &+ \sum_{i=\ell}^{s-1} \alpha_k(V_k(g(J_{k,i})) - V_k(r(J_{k,i}))). \end{aligned} \quad (9)$$

Since chunks  $J_{k,\ell+1}, \dots, J_{k,s}$  are released due to transition (8), Lemma 2 implies that  $V_k(r(J_{k,i+1})) = V_k(g(J_{k,i}))$  for all  $\ell \leq i < s - 1$ . Substituting this into (9)

$$\begin{aligned} W(A_k, r(J_{k,\ell}), d(J_{k,s})) &= W(J_{k,s}, r(J_{k,\ell}), d(J_{k,s})) \\ &+ \sum_{i=\ell}^{s-1} \alpha_k(V_k(r(J_{k,i+1})) - V_k(r(J_{k,i}))). \end{aligned} \quad (10)$$

By the telescoping summation above, it may be shown that  $W(A_k, r(J_{k,\ell}), d(J_{k,s}))$  equals  $W(J_{k,s}, r(J_{k,\ell}), d(J_{k,s})) +$

$\alpha_k(V_k(r(J_{k,s})) - V_k(r(J_{k,\ell})))$ . By the antecedent of the lemma,  $V_k(r(J_{k,\ell}))$  equals  $r(J_{k,\ell})$ . Thus

$$\begin{aligned} & W(A_k, r(J_{k,\ell}), d(J_{k,\ell})) \\ &= W(J_{k,s}, r(J_{k,\ell}), d(J_{k,s})) \\ & \quad + \alpha_k(V_k(r(J_{k,s})) - r(J_{k,\ell})). \end{aligned} \quad (11)$$

It remains to determine  $W(J_{k,s}, r(J_{k,\ell}), d(J_{k,s}))$  which is dependent on whether  $J_{k,s}$  meets its deadline. If  $J_{k,s}$  meets its deadline, then  $W(J_{k,s}, r(J_{k,\ell}), d(J_{k,s}))$  equals  $\alpha_k(V_k(g(J_{k,s})) - V_k(r(J_{k,s})))$ . Lemma 1 implies that  $V_k(g(J_{k,s})) \leq d(J_{k,s})$ ; so,  $W(J_{k,s}, r(J_{k,\ell}), d(J_{k,s}))$  does not exceed  $\alpha_k(d(J_{k,s}) - V_k(r(J_{k,s})))$ . Combining this fact and (11) implies (8) of the lemma. Therefore, the lemma is satisfied when  $J_{k,s}$  meets its deadline.

Now, consider the case where  $J_{k,s}$  misses its deadline. By Definition 1,  $W(J_{k,s}, r(J_{k,\ell}), d(J_{k,s}))$  equals  $\alpha_k P_k$ . Observe that by antecedent of the lemma,  $J_{k,s}$  is released due to transition (8); either rule (vi), (vii), or (ix) will be used to set  $d(J_{k,s})$ . Each of these rules sets  $d(J_{k,s}) = V_k(g(J_{k,s-1})) + P_k \Rightarrow P_k = d(J_{k,s}) - V_k(g(J_{k,s-1}))$ . Substituting the value of  $P_k$  and observing by Lemma 2 that  $V_k(g(J_{k,s-1}))$  equals  $V_k(r(J_{k,s}))$ , we derive  $W(J_{k,s}, r(J_{k,\ell}), d(J_{k,s}))$  equals  $\alpha_k(d(J_{k,s}) - V_k(g(J_{k,s-1})))$ . Finally, substituting the new expression for  $W(J_{k,s}, r(J_{k,\ell}), d(J_{k,s}))$  into (11) and canceling terms gives us (8) of the lemma. Thus, the lemma is also satisfied when  $J_{k,s}$  misses its deadline. ■

### E. Admission Control Algorithms

Adapting the proofs from the EDF+SRP feasibility tests in [6] and [3] for the case where application chunks, instead of jobs, are the items to be scheduled, we can find a direct mapping relation between resource-holding-times of applications and critical section lengths of jobs. The maximum blocking experienced by  $J_{k,\ell}$  is then

$$B_k = \max_{P_j > P_k} \{H_j(R_\ell) \mid \exists H_x(R_\ell) \neq 0 \wedge P_x \leq P_k\}. \quad (12)$$

In other words, the maximum amount of time for which  $J_{k,\ell}$  can be blocked is equal to the maximum resource-holding-time among all applications having a server period  $> P_k$  and sharing a global resource with some application having a server period  $\leq P_k$ . The following test may be used when the admission control algorithm has information from each application  $A_k$  on which global resources  $R_\ell$  are accessed and what the value of  $H_k(R_\ell)$  is.

*Theorem 2:* Applications  $A_1, \dots, A_q$  may be composed upon a unit-capacity processor together without any server missing a deadline, if

$$\forall k \in \{1, \dots, q\} : \sum_{P_i \leq P_k} \alpha_i + \frac{B_k}{P_k} \leq 1 \quad (13)$$

where the blocking term  $B_k$  is defined in (12).

*Proof:* This test is similar to the EDF+SRP feasibility tests in [6] and [3], substituting jobs and critical section lengths with, respectively, application chunks and resource-holding-times. We prove the contrapositive of the theorem.

Assume that the first deadline miss for some server chunk occurs at time  $t_{\text{miss}}$ . Let  $t'$  be the latest time prior to  $t_{\text{miss}}$  such that

there is no executing (and, therefore, no contending) application with deadline before  $t_{\text{miss}}$ ; since there exists an executing server from  $t'$  to  $t_{\text{miss}}$ , the processor is continuously busy in the interval  $[t', t_{\text{miss}}]$ . Observe that  $t'$  is guaranteed to exist at system-start time. The total demand imposed by server chunks in  $[t', t_{\text{miss}}]$  is defined as the sum of the execution costs of all chunks entirely contained in that interval, i.e.,  $\sum_{i=1}^q W(A_i, t', t_{\text{miss}})$ .

We will now show that the demand of any application  $A_k$  does not exceed  $\alpha_k(t_{\text{miss}} - t')$ . Let  $Y \stackrel{\text{def}}{=} \{J_{k,\ell}, \dots, J_{k,s}\}$  be the set of server chunks that the server for  $A_k$  releases in the interval  $[t', t_{\text{miss}}]$  with deadlines prior or equal to  $t_{\text{miss}}$ . If  $Y$  is empty, then the demand trivially does not exceed  $\alpha_k(t_{\text{miss}} - t')$ ; so, assume that  $Y$  is nonempty. Since the server for  $A_k$  is not in the *Executing* (or *Contending*) state immediately prior  $t'$ , it is either in the *Noncontending*, *Inactive*, or *Suspended* state for a nonzero-length time interval prior to  $t'$  [note this disallows the instantaneous transition of (6) and (8)]; therefore, the first chunk of  $Y$  must have been generated due to either transition (1) or (8), in which case either rule (i) or rule (viii) apply. Thus  $J_{k,\ell}$  is the first chunk in  $Y$ ,  $V_k(r(J_{k,\ell})) = r(J_{k,\ell})$ . We may thus partition  $Y$  into  $p$  disjoint subsequences of *successively generated* chunks  $Y^{(1)}, Y^{(2)}, \dots, Y^{(p)}$ , where  $Y^{(i)} \stackrel{\text{def}}{=} \{J_{k,\ell_i}^{(i)}, \dots, J_{k,s_i}^{(i)}\}$ . For each  $Y^{(i)}$ ,  $J_{k,\ell_i}^{(i)}$  has  $V_k(r(J_{k,\ell_i}^{(i)}))$  equal to  $r(J_{k,\ell_i}^{(i)})$ , and  $J_{k,\ell_i+1}^{(i)}, \dots, J_{k,s_i}^{(i)}$  are all released due to transition (8). Observe the chunks of each subsequence  $Y^{(i)}$  span the interval  $[r(J_{k,\ell_i}^{(i)}), d(J_{k,s_i}^{(i)})] \subseteq [t', t_{\text{miss}}]$ . By Lemma 3, the demand of  $A_k$  over the subinterval  $[r(J_{k,\ell_i}^{(i)}), d(J_{k,s_i}^{(i)})]$  does not exceed  $\alpha_k(d(J_{k,s_i}^{(i)}) - r(J_{k,\ell_i}^{(i)}))$ . Furthermore, the server for  $A_k$  does not execute in intervals not covered by the chunks of some subsequence  $Y^{(i)}$ . Since  $Y^{(1)}, Y^{(2)}, \dots, Y^{(p)}$  is a partition of  $Y$ , the subintervals do not overlap; this implies that  $\sum_{i=1}^p (d(J_{k,s_i}^{(i)}) - r(J_{k,\ell_i}^{(i)})) \leq (t_{\text{miss}} - t')$ . Therefore the total demand of  $A_k$  over interval  $[t', t_{\text{miss}}]$  does not exceed  $\alpha_k[\sum_{i=1}^p (d(J_{k,s_i}^{(i)}) - r(J_{k,\ell_i}^{(i)}))] \leq \alpha_k(t_{\text{miss}} - t')$ .

Notice that only applications with period less than  $(t_{\text{miss}} - t')$  can release chunks inside the interval: since any application  $A_k$  is not backlogged at time  $t'$ , the first chunk released after  $t'$  will have deadline at least  $t' + P_k$ . Let  $A_k$  be the application with the largest  $P_k \leq (t_{\text{miss}} - t')$ . For Theorem 1, at most one server chunk with deadline later than  $t_{\text{miss}}$  can execute in the considered interval. Therefore, only one application with period larger than  $P_k$  can execute, for at most the length of one resource-holding-time, inside the interval. The maximum amount of time that an application with period larger than  $P_k$  can execute in  $[t', t_{\text{miss}}]$  is quantified by  $B_k$ .

Since some server chunk missed a deadline at time  $t_{\text{miss}}$ , the demand in  $[t', t_{\text{miss}}]$ , plus the blocking term  $B_k$  as defined in (12), must exceed  $(t_{\text{miss}} - t')$

$$\sum_{P_i \leq P_k} (t_{\text{miss}} - t') \alpha_i + B_k \geq (t_{\text{miss}} - t'). \quad (14)$$

Dividing by  $(t_{\text{miss}} - t')$ , and then observing that  $(t_{\text{miss}} - t') \geq P_k$ , we have

$$\sum_{P_i \leq P_k} \alpha_k + \frac{B_k}{P_k} \geq 1 \quad (15)$$

which contradicts (13). ■

However, such an exact admission control test based on a policy of considering all resource usages (as the theorem above) has drawbacks. One reason is that it requires the system to keep track of each application's resource-hold times. An even more serious drawback of the more exact approach is how to fairly account for the "cost" of admitting an application into the open environment. For example, an application that needs a VP speed twice that of another should be considered to have a greater cost (all other things being equal); considered in economic terms, the first application should be "charged" more than the second, since it is using a greater fraction of the platform resources and thus having a greater (adverse) impact on the platform's ability to admit other applications at a later point in time.

However, in order to measure the impact of global resource-sharing on platform resources, we need to consider the resource usage of not just an application, but of all other applications in the systems. Consider the following scenario. If application  $A_1$  is using a global resource that no other application chooses to use, then this resource usage has no adverse impact on the platform. Now, if a new application  $A_2$  with a very small period parameter that needs this resource seeks admission, the impact of  $A_1$ 's resource-usage becomes extremely significant (since  $A_1$  would, according to the SRP, block  $A_2$  and also all other applications that have a period parameter between  $A_1$ 's and  $A_2$ 's). So, how should we determine the cost of  $A_1$ 's use of this resource, particularly if we do not know beforehand whether or not  $A_2$  will request admission at a later point in time?

To sidestep the dilemma described above, we believe a good design choice is to effectively *ignore* the exact resource-usage of the applications in the online setting, instead considering only the maximum amount of time for which an application may choose to hold any resource; also, we did not consider the identity of this resource. That is, we required a simpler interface than the one discussed in Section II, in that rather than requiring each application to reveal its maximum resource-holding times on all  $m$  resources, we only require each application  $A_k$  to specify a *single* resource-holding parameter  $H_k$ , which is defined as follows:

$$H_k \stackrel{\text{def}}{=} \max_{\ell=1}^m H_k(R_\ell). \quad (16)$$

The interpretation is that  $A_k$  may hold *any* global resource for up to  $H_k$  units of execution. With such characterization of each application's usage of global resources, we ensure that we do not admit an application that would unfairly block other applications from executing due its large resource usage. This test, too, is derived directly from the EDF+SRP feasibility test of Theorem 2, and is as follows:

---

#### Algorithm Admit ( $A_k = (\alpha_k, P_k, H_k)$ )

---

- 1) **for** each  $A_i: (P_i \geq P_k)$  **do**
- 2)   **if**  $\max_{j: P_j > P_i} H_j > P_i(1 - \sum_{j: P_j \leq P_i} \alpha_j)$   
      **return** "reject"
- 3) **for** each  $A_i: (P_i < P_k)$  **do**
- 4)   **if**  $H_k > P_i(1 - \sum_{P_j \leq P_i} \alpha_j)$   
      **return** "reject"
- 5) **return** "admit"

It follows from the properties of the SRP, (as proved in [3]) that the new application  $A_k$ , if admitted, may *block* the execution of applications  $A_i$  with period parameter  $P_i < P_k$ . Moreover, by Theorem 2, it may *interfere* with applications  $A_i$  with period parameter  $P_i \geq P_k$ . Since the maximum amount by which any application  $A_j$  with  $P_j > P_i$  may block an application  $A_i$  is equal to  $H_j$ , lines 1–2 of Procedure ALGORITHM ADMIT determine whether this blocking can cause any application  $A_i$  with  $P_i \geq P_k$  to miss its deadline. Similarly, since the maximum amount by which application  $A_k$  may block any other application is, by definition of the interface, equal to  $H_k$ , lines 3–4 of ALGORITHM ADMIT determine whether  $A_k$ 's blocking causes any other application with  $P_i < P_k$  to miss its deadline. If the answer in both cases is "no," then ALGORITHM ADMIT admits application  $A_k$  in line 5.

#### F. Enforcement

One of the major goals in designing open environments is to provide interapplication *isolation*—all other applications should remain unaffected by the behavior of a misbehaving application. By encapsulating each application into a BROE server, we provide the required isolation, enforcing a correct behavior for every application.

Using techniques similar to those used to prove isolation properties in CBS-like environments (see, e.g., [1] and [27]), it can be shown that our open environment does indeed guarantee interapplication isolation in the absence of resource-sharing. It remains to study the effect of resource-sharing on interapplication isolation.

Clearly, applications that share certain kinds of resources cannot be completely isolated from each other: for example if one application corrupts a shared data-structure then all the applications sharing that data structure are affected. When a resource is left in an inconsistent state, one option could be to inflate the resource-holding time parameters with the time needed to reset the shared object to a consistent state, when there is such a possibility.

However, we believe that it is rare that truly independently developed applications share "corruptible" objects—good programming practice dictates that independently developed applications do not depend upon proper behavior of other applications (and in fact this is often enforced by operating systems). Hence the kinds of resources we expect to see shared between different applications are those that the individual applications cannot corrupt. In that case, the only misbehavior of an application  $A_k$  that may affect other applications is if it holds on to a global resource for greater than  $\alpha_k P_k$ , or than the  $H_k$  time units of execution that it had specified in its interface. To prevent this, we assume that our enforcement algorithm simply preempts  $A_k$  after it has held a global resource for  $\min\{H_k, \alpha_k P_k\}$ , and ejects it from the shared resource. This may result in  $A_k$ 's internal state getting compromised, but the rest of the applications are not affected. Note that such an enforcement algorithm might require to set a timer each time a resource is locked, increasing the system overhead. We believe this is the minimum price to pay for guaranteeing temporal isolation among applications that share global resources. Applications that repeatedly access global resources will need to be charged with a larger

overhead, unless performing multiple accesses within a single lock.

When applications do share corruptible resources, we have argued above that isolation is not an achievable goal; however, *containment* [19] is. The objective in containment is to ensure that the only applications affected by a misbehaving application are those that share corruptible global resources with it—the intuition is that such applications are not truly independent of each other. We have strategies for achieving some degree of containment; for instance, one option could be to donate to an application locking a corruptible resource the bandwidth of applications that share the same resource. However, discussion of these strategies is beyond the scope of this document.

### G. Bounded Delay Property

The bounded-delay resource partition model, introduced by Mok *et al.* [35], is an abstraction that quantifies resource “supply” that an application receives from a given resource.

*Definition 3:* A server implements a *bounded-delay* partition  $(\alpha_k, \Delta_k)$  if in any time interval of length  $L$  during which the server is continually backlogged, it receives at least

$$(L - \Delta_k)\alpha_k$$

units of execution.

*Definition 4:* A *bounded-delay server* is a server that implements a bounded-delay partition.

We will show that when every application is admitted through a proper admission control test, BROE implements a bounded delay partition. Before proving this property, we need some intermediate lemmas. The first lemma quantifies the minimum virtual-time  $V_k$  for a server for application  $A_k$  that is in the *Contending* or *Executing* state.

*Lemma 4:* Given BROE servers of applications  $A_1, \dots, A_q$  satisfying Theorem 2, if server chunk  $J_{k,\ell}$  of server  $A_k$  is executing or contending at time  $t$  (where  $r(J_{k,\ell}) \leq t \leq d(J_{k,\ell})$ ), then

$$V_k(t) \geq V_k(r(J_{k,\ell})) + \frac{\max(0, t - V_k(r(J_{k,\ell})) - P_k(1 - \alpha_k))}{\alpha_k}. \quad (17)$$

*Proof:* The proof is by contradiction. Assume that all servers have been admitted to the open environment via Theorem 2, but there exists a server  $A_k$  in the *Contending* or *Executing* state at time  $t$  that has

$$V_k(t) < V_k(r(J_{k,\ell})) + \frac{\max(0, t - V_k(r(J_{k,\ell})) - P_k(1 - \alpha_k))}{\alpha_k}. \quad (18)$$

Since  $V_k$  never decreases, the above strict inequality implies that

$$t > V_k(r(J_{k,\ell})) + P_k(1 - \alpha_k). \quad (19)$$

We will show that if (18) holds, there exist a legal scenario under which  $A_k$  will miss a server deadline. Assume that application  $A_k$  has  $\alpha_k P_k$  units of execution backlogged at time  $r(J_{k,\ell})$  (the server can be in any state immediately prior to

$r(J_{k,\ell})$ ); also assume that no job of application  $A_k$  requests any global resources during the next  $\alpha_k P_k$  units of  $A_k$ 's execution (i.e., transition (6) will not be used). The described scenario is a legal scenario for application  $A_k$  with parameter  $\alpha_k$  and  $\Delta_k$ .

Note that each of the server deadline update rules essentially sets  $D_k$  equal to  $V_k + P_k$ ; therefore,  $d(J_{k,\ell})$  equals  $V_k(r(J_{k,\ell})) + P_k$ . The current time remaining until  $J_{k,\ell}$ 's deadline is  $V_k(r(J_{k,\ell})) + P_k - t$ . The virtual time of  $A_k$  at time  $t$  by (18) and (19) satisfies the following inequality:

$$V_k(t) < V_k(r(J_{k,\ell})) + \frac{1}{\alpha_k} \cdot (t - V_k(r(J_{k,\ell})) - P_k(1 - \alpha_k)). \quad (20)$$

The remaining amount of time at time  $t$  that the server for  $A_k$  must execute for  $V_k$  to equal  $d(J_{k,\ell})$  (i.e., complete its execution) is  $\alpha_k(V_k(r(J_{k,\ell})) + P_k - V_k(t))$ . Combining this expression with (20), the remaining execution time is strictly greater than  $V_k(r(J_{k,\ell})) + P_k - t$ . However, this exceeds the remaining time to the deadline; since the server for  $A_k$  is continuously in the *Contending* or *Executing* state throughout this scenario, the server will miss a deadline at  $d(J_{k,\ell})$ . This contradicts the lemma, given that the servers satisfied Theorem 2. Our original supposition of (18) is falsified and the lemma follows. ■

We next show that at any time a server chunk is released for  $A_k$ , the actual time must exceed the virtual time.

*Lemma 5:* For any server chunk  $J_{k,\ell}$  of the BROE server for  $A_k$ ,

$$V_k(r(J_{k,\ell})) - r(J_{k,\ell}) \leq 0. \quad (21)$$

*Proof:* The lemma may be proved by analyzing each of the server rules involved in moving the server state to *Contending*. If the server state for  $A_k$  is inactive prior to the release of chunk  $J_{k,\ell}$ , then rule (i) sets  $V_k$  to current time, and the lemma is satisfied. If the server was suspended immediately prior to the release of  $J_{k,\ell}$ , rule (viii) releases  $J_{k,\ell}$  only when  $Z_k$  equals  $t_{\text{cur}}$ . Observe that all the rules of the server set  $Z_k$  to a value greater than or equal to  $V_k$ . Thus,  $V_k(r(J_{k,\ell})) - r(J_{k,\ell})$  is either zero or negative. ■

The final lemma before proving that BROE is a bound-delay server, shows that for any server the absolute difference between virtual time and actual time is bounded in terms of the server parameters.

*Lemma 6:* For an application  $A_k$  admitted in the open environment, if the server for  $A_k$  is backlogged at time  $t$ , then

$$|V_k(t) - t| \leq P_k(1 - \alpha_k). \quad (22)$$

*Proof:* If the server for  $A_k$  is in the *Suspended* state, then because the server is backlogged this implies that  $V_k(t) - t > 0$ ; so, the server will not become contending until time  $V_k$ . Let  $t'$  be the last time prior to  $t$  that the server was contending or executing; it's easy to see that  $V_k(t') - t' > V_k(t) - t$ . So, we will reason about  $t'$  and show for any such contending or executing time  $V_k(t') - t' \leq P_k(1 - \alpha_k)$ . If the server for  $A_k$  was contending or executing at time  $t$ , let  $t'$  instead equal  $t$ . We will show in the remainder of the proof that  $-P_k(1 - \alpha_k) \leq V(t') - t' \leq P_k(1 - \alpha_k)$ . Let  $J_{k,\ell}$  be the server chunk corresponding to the last contending or executing state at  $t'$  for application  $A_k$ . Observe that this implies that  $t' \geq r(J_{k,\ell})$ .

Let us first show that  $V(t') - t' \leq P_k(1 - \alpha_k)$ . Observe that because virtual time progresses at a rate equal to  $1/\alpha_k$  and cannot exceed  $D_k = V_k(r(J_{k,\ell})) + P_k = V_k(r(J_{k,\ell})) + 1/\alpha_k \cdot (\alpha_k P_k)$

$$V_k(t') \leq V_k(r(J_{k,\ell})) + \frac{1}{\alpha_k} \cdot \min(t' - r(J_{k,\ell}), \alpha_k P_k). \quad (23)$$

Subtracting  $t'$  from both sides, and observing that the RHS is maximized for  $t' = r(J_{k,\ell}) + \alpha_k P_k$  (since  $t' \geq r(J_{k,\ell})$ , the “min” term increases at a rate of  $1/\alpha_k$  in  $[r(J_{k,\ell}), r(J_{k,\ell}) + \alpha_k P_k]$ , and remains constant for all  $t > r(J_{k,\ell}) + \alpha_k P_k$ ) implies

$$V_k(t') - t' \leq V_k(r(J_{k,\ell})) + P_k - r(J_{k,\ell}) - \alpha_k P_k. \quad (24)$$

Lemma 5 implies that  $V_k(r(J_{k,\ell})) - r(J_{k,\ell}) \leq 0$ . Thus, (24) may be written as  $V_k(t') - t' \leq P_k - \alpha_k P_k = P_k(1 - \alpha_k)$ , proving the upper bound on  $V_k(t') - t'$  (and thus an upper bound on  $V(t) - t$ ).

We will now prove a lower bound on  $V(t') - t'$ . By Lemma 4 and the fact that the server is contending or executing at time  $t'$ , we have a lower bound on the virtual time at  $t'$

$$V_k(t') - t' \geq V_k(r(J_{k,\ell})) + \frac{1}{\alpha_k} \cdot \max(0, t' - V_k(r(J_{k,\ell})) - P_k(1 - \alpha_k)) - t'. \quad (25)$$

Observe that for all  $t' : r(J_{k,\ell}) \leq t' \leq V_k(r(J_{k,\ell})) + P_k(1 - \alpha_k)$ , the “max” term in the RHS of the above expression is zero; thus, the RHS decreases in  $t'$  from  $r(J_{k,\ell})$  to  $V_k(r(J_{k,\ell})) + P_k(1 - \alpha_k)$ . For  $t'$  greater than  $V_k(r(J_{k,\ell})) + P_k(1 - \alpha_k)$ , the “max” term increases at a rate of  $(1)/(\alpha_k)$ . Therefore, the RHS of the above inequality is minimized when  $t'$  equals  $V_k(r(J_{k,\ell})) + P_k(1 - \alpha_k)$ . Thus,  $V_k(t) - t' \geq -P_k(1 - \alpha_k)$ , proving the lower bound on  $V(t) - t$ ; the lemma follows. ■

We are now ready to prove that BROE implements a bounded-delay partition.

*Theorem 3 (Bounded-Delay Property):* BROE is a bounded-delay server.

*Proof:* From the definition of the BROE server, it can be seen that the virtual time  $V_k$  is updated only when an inactive application goes active or whenever subsequently it is executing. In the latter case,  $V_k$  is incremented at a  $1/\alpha_k$  rate. Thus, there is a direct relation between the execution time allocated to the application through the BROE server and the supply the application would have received if scheduled on a VP of speed  $\alpha_k$ . The quantity  $V_k(t) - t$  indicates the advantage the application  $A_k$  executing through the BROE server has compared with VP in terms of supply. If the above term is positive, the application received more execution time than the VP would have by time  $t$ . If it is negative, the BROE server is “late.”

From Lemma 6, the execution time supplied to an application through a dedicated BROE server never exceeds nor is exceeded by the execution time it would have received on a dedicated VP for more than  $P_k(1 - \alpha_k)$  time units. The “worst case” is when both displacements happen together, i.e., interval  $T$  starts when  $V_k(t) - t = P_k(1 - \alpha_k)$  and ends when  $V_k(t) - t = -P_k(1 - \alpha_k)$ . This interval in which the BROE server can be delayed from executing, while still satisfying the bound on  $|V(t) - t|$  from

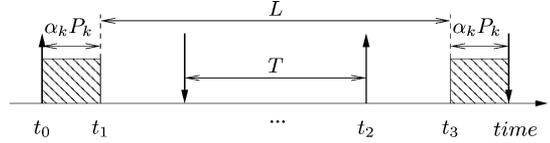


Fig. 3. Worst case scenario discussed in the Proof of Theorem 3. The application receives execution during the shaded intervals.

Lemma 6 is of length at most twice  $P_k(1 - \alpha_k)$ . By the definition of  $P_k$  (1), this is equal to  $\Delta_k$ . Thus, the maximum delay that an application executing on a BROE server may experience is  $\Delta_k$ .

In other words, it can be shown that the “worst case” (see Fig. 3) occurs when application  $A_k$ :

- receives execution immediately upon entering the *Contending* state (at time  $t_0$  in the figure), and the interval of length  $L$  begins when it completes execution and undertakes transition (6) to the *Suspended* state (at time  $t_1$  in the figure);
- after having transited between the *Suspended*, *Contending*, and *Executing* states an arbitrary number of times, undertakes transition (8) to enter the *Contending* state (time  $t_2$  in the figure) at which time it is scheduled for execution [transition (2)] as late as possible; the interval ends just prior to  $A_k$  being selected for execution (time  $t_3$  in the figure).

A job arriving at time  $t_1$  will be served by the BROE with the maximum delay of  $\Delta_k$  from the supply granted by a VP of speed  $\alpha_k$ . Since the execution received in an interval  $T$  going from the deadline of the first chunk (released at  $t_0$ ) until the release time of the last one at  $t_2$  cannot be higher than  $\alpha_k T$ , the supply granted over interval  $L$  is  $\alpha_k T = (L - 2P_k(1 - \alpha_k))\alpha_k$ . By the definition of  $P_k$  (1), this is equal to  $(L - \Delta_k)\alpha_k$ . ■

#### IV. APPLICATION-LEVEL SCHEDULERS

In the previous section, we analyzed how to compose multiple servers on the same processor without violating the bounded-delay server constraints. Provided that these *global* constraints are met, we now address the *local* schedulability problem, to verify if a collection of jobs composing an application can be scheduled on a bounded delay server with given  $\alpha_k$  and  $\Delta_k$ , when jobs can share exclusive resources with other applications.

To do this, we have three options on how to schedule and validate the considered collection of jobs.

- 1) Validate the application on a dedicated VP with speed  $\alpha_k$  using a given scheduling algorithm. If every job is completed at least  $\Delta_k$  time-units before its deadline, then the application is schedulable on a bounded-delay partition  $(\alpha_k, \Delta_k)$  when jobs are scheduled according to the same order as they would on a dedicated VP schedule.
- 2) Validate the application on a dedicated VP with speed  $\alpha_k$  using EDF. If every job is completed at least  $\Delta_k$  time-units before its deadline, then the application is schedulable with EDF on a bounded-delay partition  $(\alpha_k, \Delta_k)$ , without needing to “copy” the VP schedule.
- 3) Validate the application by analyzing the execution time effectively supplied by the partition in the worst-case and the demand imposed by the jobs scheduled with any scheduling algorithm, avoiding validation on a VP.

These options are hereafter explained in more detail.

### A. Replicating the VP Scheduling

When scheduling a set of applications on a shared processor, there is sometimes the need to preserve the original scheduling algorithm with which an application has been conceived and validated on a slower processor. If this is the case, we need to guarantee that all jobs composing the application will still be schedulable on the bounded-delay partition provided by the open environment through the associated BROE server. Mok *et al.* [22], [35] have previously addressed this problem. We restate their result, adapting it to the notation used so far.

*Theorem 4: ([35, Th. 6]):* Given an application  $A_k$  and a Bounded Delay Partition  $(\alpha_k, \Delta_k)$ , let  $S_n$  denote a valid schedule on a VP with speed  $\alpha_k$ , and  $S_p$  the schedule of  $A_k$  on partition  $(\alpha_k, \Delta_k)$  according to the same execution order and amount as  $S_n$ . Also let  $\bar{\Delta}_k$  denote the largest amount of time such that any job of  $S_n$  is completed at least  $\bar{\Delta}_k$  time units before its deadline.  $S_p$  is a valid schedule if and only if  $\bar{\Delta}_k \geq \Delta_k$ .

The theorem states that all jobs composing an application are schedulable on a BROE server having  $\alpha_k$  equal to the VP speed and  $\Delta_k$  equal to the jitter tolerance of the VP schedule, provided that jobs are executed in the *same execution order* of the VP schedule.

In order to be applicable to general systems, this approach would require that each individual application's scheduling event (job arrivals and completions) be "buffered" during the delay bound  $\Delta_k$ —essentially, an event at time  $t_o$  is ignored until the earliest time-instant when  $V(t) \geq t$ —so that events are processed in the same order in the open environment as they would be if each application were running upon its dedicated VP. However, we will see that such buffering is unnecessary when the individual application can be EDF-scheduled in the open environment.

### B. Application-Level Scheduling Using EDF

To avoid the complexity of using buffers to keep track of the scheduling events, it is possible to use a simplified approach. When an application does not mandate to be scheduled with a particular scheduling algorithm, we show that EDF can be optimally used as application-level scheduler for the partition, without needing to "copy" the VP behavior. To distinguish the buffered from the native version of the partition local scheduler, we will call VP-EDF the application-level scheduler reproducing the VP behavior, while the normal local scheduler using only jobs earliest deadlines will be simply called EDF.

*Definition 5:* A scheduling algorithm is *resource-burst-robust* if advancing the supply, the schedulability is preserved.

*Lemma 7 (From Feng [21]):* EDF is resource-burst-robust.

*Lemma 8:* Consider an application  $A_k$ , composed by a set of jobs with fixed release times and deadlines. If all jobs of application  $A_k$  always complete execution at least  $\Delta_k$  time units prior to their deadlines when scheduled with EDF upon a dedicated VP of computing capacity  $\alpha_k$ , then all jobs of  $A_k$  are schedulable with EDF on a partition  $(\alpha_k, \Delta_k)$ .

*Proof:* We prove the contrapositive. Assume a collection of jobs of an application  $A_k$  completes execution at least  $\Delta_k$  time units prior to their deadlines when scheduled with EDF on a

dedicated  $\alpha_k$ -speed VP, but some of these jobs miss a deadline when  $A_k$  is scheduled with EDF on a partition  $(\alpha_k, \Delta_k)$ . Let  $t_{\text{miss}}$  be the first time a deadline is missed and let  $t_s$  denote the latest time-instant prior to  $t_{\text{miss}}$  at which there are no jobs with deadline  $\leq t_{\text{miss}}$  awaiting execution in the partition schedule ( $t_s \leftarrow 0$  if there was no such instant). Hence, over  $[t_s, t_{\text{miss}})$ , the partition is only executing jobs with deadline  $\leq t_{\text{miss}}$ , or jobs that were *blocking* the execution of jobs with deadline  $\leq t_{\text{miss}}$ . Let  $Y$  be the set of such jobs.

Since a deadline is missed, the total amount of demand of jobs in  $Y$  during  $[t_s, t_{\text{miss}})$  upon the BROE server is greater than the execution time supplied in the same interval. From Lemma 7, we know that the minimum amount of execution  $A_k$  would receive in interval  $[t_s, t_{\text{miss}})$ , is  $\alpha_k((t_{\text{miss}} - t_s) - \Delta_k)$ .

Consider now the VP schedule. Since every job completes at least  $\Delta_k$  time-units before its deadline, the job that misses its deadline in the partition schedule will complete before instant  $t_{\text{miss}} - \Delta_k$  in the VP schedule. Moreover, since EDF always schedules tasks according to their absolute deadline, no jobs in  $Y$  will be scheduled in interval  $[t_{\text{miss}} - \Delta_k, t_{\text{miss}}]$ . Therefore, the total demand of jobs in  $Y$  during  $[t_s, t_{\text{miss}})$  does not exceed  $\alpha_k((t_{\text{miss}} - \Delta_k) - t_s)$ . However, this contradicts the fact that the minimum amount of execution that is provided by the BROE server over this interval is  $\alpha_k((t_{\text{miss}} - \Delta_k) - t_s)$ . ■

Lemma 8 is a stronger result than the one in [35, Corollary 4], where applications needed to be scheduled according to VP-EDF. In [21, Th. 2.7], a more general result is proved, saying that any resource-burst-robust scheduler can be used without needing to reproduce the VP schedule. However, there is a flaw in this result; for instance, even though DM is a resource-burst-robust scheduler, it cannot be used without buffering events. To see this, consider the following example.

*Example 1:* An application composed of two periodic tasks  $\tau_1 = (1, 6, 4)$  and  $\tau_2 = (1, 6, 6)$  is validated on a processor of speed  $\alpha_k = 1/2$ . Each job is completed at least  $\Delta_k = 2$  time units prior to its deadline. However, when  $A_k$  is scheduled on a bounded-delay partition  $(\alpha_k, \Delta_k) = (1/2, 2)$ , it can miss a deadline when both  $\tau_1$  and  $\tau_2$  release jobs at time  $t$  and the server temporarily exhausts its budget. The application may have to wait until time  $t + \Delta = t + 2$  to receive service, at which point  $\tau_1$  executes in  $[t + 2, t + 3)$  (exhausting the budget). The next service interval is  $[t + 4, t + 5)$ , when the next job of  $\tau_1$  is scheduled. Then,  $\tau_2$  has to wait for the next service interval  $[t + 6, t + 7)$ , but at that point it would miss its deadline. ■

Notice that since the Proof of Lemma 8 does not rely on any particular protocol for the access to shared resources, the validity of the result can be extended to every reasonable policy, like SRP [3] or others, provided that the same mechanism is used for both the VP and the partition schedule. Since EDF+SRP is an optimal scheduling algorithm for VPs [6], the next theorem follows.

*Theorem 5:* A collection of jobs is schedulable with EDF+SRP on a partition  $(\alpha_k, \Delta_k)$  if and only if it is schedulable with some scheduling algorithm on an  $\alpha_k$ -speed VP with a jitter tolerance of  $\Delta_k$ , i.e., all jobs finish at least  $\Delta_k$  time units before their deadline.

Therefore, when there is no limit on the algorithm to be used to schedule the application jobs on a partition, using EDF+SRP

is an optimal choice, since it guarantees that all deadlines are met independently from the algorithm that has been used for the validation on the dedicated VP. This also explains the meaning of the names we gave in Section II to  $\alpha_k$  and  $\Delta_k$  parameters.

On the contrary, when the scheduling algorithm cannot be freely chosen, for instance when a fixed priority order among tasks composing an application has to be enforced, we showed in Section IV-A that a buffered version of the VP schedule can be used. However, to avoid the computational effort of reproducing the VP scheduling order at runtime, some more expense can be paid offline by analyzing the execution time supplied by the partition together with the demand imposed by the jobs of the application. The next section addresses this problem.

### C. Application-Level Scheduling With Other Algorithms

The application may require that a scheduler other than EDF+SRP be used as an application-level scheduler. When a buffered version of the VP schedule is not feasible due to the associated runtime complexity, an alternative could be to use a more sophisticated schedulability analysis instead of the validation process on a dedicated VP. This requires one to consider the service effectively supplied by the open environment in relation to the amount of execution requested by the application. Our BROE server implements a bounded-delay server in the presence of shared resources. Examples of analysis for the fixed-priority case under servers implementing bounded-delay partitions or related partitions, in absence of shared resources, can be found in [28], [35], [41], [42], and easily applied to our open environment. We conjecture that the results for local fixed-priority schedulability analysis on resource partitions can be easily extended to include local and global resources, and be scheduled by BROE without modification to the server. We leave the exploration of this conjecture as a future work.

## V. LOCAL SCHEDULABILITY ANALYSIS

Since we showed (Theorem 5) that EDF+SRP can be used to optimally schedule the jobs composing an application on a bounded-delay partition, we choose this algorithm as the default local scheduling algorithm for our open environment. This means that whenever an application does not require a different scheduling policy, the open environment will schedule the jobs of an admitted executing application using EDF with resource access arbitrated through the SRP protocol, since this choice allows optimizing system performances. We will hereafter derive local schedulability tests for this case.

Since we are moving the detail of our analysis from server chunks to the jobs composing an application, we need to extend the notational model used. In the rest of this paper, we will consider an application  $A_k$  to be composed of  $n_k$  sporadic tasks [8].<sup>4</sup> We will indicate with  $c_i, d_i, p_i$ , the WCET, relative deadline, and period or minimum interarrival time of task  $\tau_i$ . The maximum size of a critical section on resource  $R_\ell$  accessed by a task  $\tau_i$  is denoted with  $h_i(R_\ell)$ . To avoid confusion with the shared resource policy adopted at system-level and described in Section III-A, we will distinguish between *local* and *global*

<sup>4</sup>We believe that the following analysis can be easily adapted also for more general models, such as arbitrary collections of jobs.

ceiling of a shared resource  $R_\ell$ . Global ceiling  $\Pi(R_\ell)$  is given by the minimum value from among all the period parameters  $P_k$  of applications  $A_k$  that use resource  $R_\ell$ ; local ceiling  $\pi_k(R_\ell)$  is given, locally to each application  $A_k$ , by the minimum value from among all relative deadlines  $d_i$  of tasks  $\tau_i \in A_k$  that can lock resource  $R_\ell$ . Note that the *system ceiling* used for global SRP—equal to the minimum global ceiling of any resource that is locked at a given instant—corresponds at the local level with the *application ceiling* for local SRP—defined as the minimum local ceiling of any resource that is locked at a given instant by tasks of the considered application.

For any collection of jobs released by an application  $A_k$  and any real number  $t \geq 0$ , the *demand bound function*  $\text{DBF}(A_k, t)$  is defined as the largest cumulative execution requirement of all jobs that can be generated by  $A_k$  to have both their arrival times and their deadlines within a contiguous interval of length  $t$ . For instance, for the sporadic task model it has been shown [8] that the cumulative execution requirement of jobs over an interval  $[t_o, t_o + t)$  is maximized if all tasks arrive at the start of the interval—i.e., at time-instant  $t_o$ —and subsequent jobs arrive as rapidly as permitted—i.e., at instants  $t_o + p_i, t_o + 2p_i, t_o + 3p_i, \dots$

$$\text{DBF}(A_k, t) \stackrel{\text{def}}{=} \sum_{\tau_i \in A_k} \max \left( 0, \left( \left\lfloor \frac{t - d_i}{p_i} \right\rfloor + 1 \right) c_i \right). \quad (26)$$

Algorithms are described in [8] for efficiently computing  $\text{DBF}(A_k, t)$  in constant time for the sporadic task model, for any  $t \geq 0$ .

Shin and Lee [42] showed that  $\text{DBF}(A_k, t)$ , together with a bound on the provided supply, can be used to derive a schedulability condition for periodic task systems scheduled on a bounded-delay resource partition  $(\alpha_k, \Delta_k)$ . We restate their result in terms of general collection of jobs released by an application  $A_k$ , whose execution requirements are captured through a demand bound function  $\text{DBF}(A_k, t)$ .

*Theorem 6 (Adapted From Shin and Lee [42]):* An application  $A_k$  is EDF-schedulable on a bounded-delay partition  $(\alpha_k, \Delta_k)$  if for all  $t \geq 0$

$$\text{DBF}(A_k, t) \leq \alpha_k(t - \Delta_k). \quad (27)$$

Baruah proposed in [6] a technique for analyzing systems scheduled with EDF+SRP using the demand-bound function. The approach is to calculate, for all  $L > 0$ , the maximum amount of time that a job with relative deadline at most  $L$  could be “blocked” by a job with relative deadline greater than  $L$ . The following function [6] quantifies this maximum blocking over an interval of length  $L$  in an application  $A_k$ :

$$b_k(L) \stackrel{\text{def}}{=} \max \{ h_i(R_j) \mid \exists \ell : (\tau_i, \tau_\ell \in A_k) \wedge (d_i > L) \wedge (d_\ell \leq L) \wedge (\tau_i \text{ and } \tau_\ell \text{ access } R_j) \}. \quad (28)$$

The next corollary follows from the application of [6, Th. 1] and Theorem 6 above.

*Corollary 1:* An application  $A_k$  is EDF+SRP-schedulable on a bounded-delay partition  $(\alpha_k, \Delta_k)$  if for all  $L > 0$

$$b_k(L) + \text{DBF}(A_k, L) \leq \alpha_k(L - \Delta_k). \quad (29)$$



of the *real time* elapsed between the lock and release of the resource.<sup>5</sup>

### B. Decreasing Resource Holding Times

Since the resource holding time determines the time for which other servers below the global ceiling of the locked resource can be blocked, it is very important to keep this value as small as possible. In this way, it is possible to compose more servers on the same system without having to account for a large server blocking time. One way to do that, is using techniques from [13], [23] and [12]. Basically, these techniques artificially decrease the ceiling of a resource by adding a dummy critical section to tasks with a deadline lower than the resource ceiling, when this does not affect the schedulability of the application. The semantics of the application do not change, but the resulting resource holding time decreases due to the reduced number of preemptions on a task holding the resource. The procedure to decrease the ceiling to the minimum possible value depends on the scheduling algorithm used. The EDF case is treated in [13], and [23], while [12] deals with RM/DM. Both works assume a dedicated computing resource is used. We need to adapt them to the case in which a task set is to be scheduled on a bounded delay server. We will consider here only the EDF case.

When the application is locally scheduled with EDF, the only modification needed to the algorithm presented in [23] is in the schedulability test used to verify if the ceiling can be decreased. Instead of the test used at line 2 of Procedure REDUCECEILING in [23], we will use the test given by Corollary 1. Since this is just a trivial modification, we do not include here a detailed description of the algorithm, which can be found in [13] and [23].

In the following section, we will provide an alternative and efficient way to decrease the resource holding time of an application, without requiring any additional computation.

### C. Executing Global Critical Sections Without Local Preemptions

In this section, we will formally show that even if application  $A_k$  was validated upon a dedicated VP of speed  $\alpha_k$  using EDF+SRP or any other policy, some critical sections may be executed without local preemptions under a BROE server.

*Theorem 7:* Given an application  $A_k$  accessing a globally shared resource  $R_\ell$  can be scheduled upon a dedicated VP of speed- $\alpha_k$ , where each job completes at least  $\Delta_k$  time units prior to its deadline: if  $h_k(R_\ell) \leq \Delta_k/2$  and  $H_k(R_\ell) \leq \alpha_k P_k$  then  $A_k$  can be EDF-scheduled on a BROE server with parameters  $(\alpha_k, \Delta_k)$ , executing each critical section of  $R_\ell$  with local preemptions disabled.

*Proof:* The proof is by contradiction. Assume the antecedent of the theorem holds, but the application  $A_k$  misses a deadline on a BROE server with parameters  $(\alpha_k, \Delta_k)$ , when executing a critical section of a global resource  $R_\ell$  with local preemptions disabled. For Theorem 5,  $A_k$  is schedulable on a partition  $(\alpha_k, \Delta_k)$  with EDF+SRP.

Let  $t_{\text{miss}}$  be the first missed deadline with nonpreemptive locking, and  $t_s$  the latest time-instant prior to  $t_{\text{miss}}$  at which

there are no jobs with deadline  $\leq t_{\text{miss}}$  awaiting execution ( $t_s \leftarrow 0$  if there was no such instant). With nonpreemptive locking, as with SRP, there is at most one blocking job over  $[t_s, t_{\text{miss}}]$  (see [6]). The blocking job must have acquired the lock before  $t_s$  and have a deadline after  $t_{\text{miss}}$ . Let  $Y$  be the set of jobs that have both release time and deadline in  $[t_s, t_{\text{miss}}]$ . Over  $[t_s, t_{\text{miss}}]$ , the application executes only jobs in  $Y$ , or a blocking job. Let  $t_l$  and  $t_u$  be, respectively, the time at which the blocking job acquires and releases the lock on  $R_\ell$ , when EDF+SRP is used as a local scheduling algorithm. We hereafter prove that  $t_{\text{miss}} \leq t_u$ . Suppose, by contradiction,  $t_{\text{miss}} > t_u$ . The total demand in  $[t_s, t_{\text{miss}}]$  is given by the contributions of the jobs in  $Y$ , plus the remaining part of the blocking critical section that has still to be executed at  $t_s$ . The demand due to jobs in  $Y$  is the same whether SRP or nonpreemptive locking is used. Moreover, since with nonpreemptive locking a job cannot be interfered with while holding a global lock, the contribution of the blocking job in  $[t_s, t_{\text{miss}}]$  cannot be larger than with SRP. Therefore, the total amount of execution requested in  $[t_s, t_{\text{miss}}]$  with SRP is at least as large as with nonpreemptive locking. Since a deadline is missed with nonpreemptive locking, a deadline will be missed even with SRP, contradicting the hypothesis. Then,  $t_{\text{miss}} \leq t_u$ .

When there is a blocking contribution in  $[t_s, t_{\text{miss}}]$ , it means that the blocking job is still holding the lock on  $R_\ell$  at time  $t_s$ . Therefore, with nonpreemptive locking, the blocking job is the only one executing in  $[t_l, t_s]$ . We now prove that no job with deadline  $\leq t_{\text{miss}}$  of  $A_k$  is released in  $[t_l, t_s]$ . Suppose that a job  $J$  is released in that interval. Since the blocking job is executing nonpreemptively in  $[t_l, t_s]$ ,  $J$  is still awaiting execution at time  $t_s$ . Since  $t_{\text{miss}}$  is the first missed deadline,  $J$ 's deadline cannot be  $< t_s$ . Moreover, for the definition of  $t_s$ ,  $J$ 's deadline cannot even be between  $t_s$  and  $t_{\text{miss}}$ . Thus,  $J$ 's deadline must be later than  $t_{\text{miss}}$ .

Let  $W(Y)$  be the total amount of execution requested in  $[t_s, t_{\text{miss}}]$  by jobs in  $Y$ , and let  $\xi^{\text{NPL}}(t_1, t_2)$  (resp.  $\xi^{\text{SRP}}(t_1, t_2)$ ) be the amount of execution received by the blocking job in  $[t_1, t_2]$  when nonpreemptive locking (resp. SRP) is used. Finally, let  $s(t_1, t_2)$  be the amount of execution supplied to  $A_k$  in interval  $[t_1, t_2]$ . Since a deadline is missed with nonpreemptive blocking, the following relation holds:

$$W(Y) + \xi^{\text{NPL}}(t_s, t_{\text{miss}}) > s(t_s, t_{\text{miss}}) = s(t_l, t_{\text{miss}}) - s(t_l, t_s). \quad (32)$$

Consider the EDF+SRP schedule. The BROE server is never suspended in  $[t_l, t_u]$  (remember  $t_u$  is the unlocking time when SRP is used), and  $A_k$  has at least one pending job throughout the same interval. Rule (ix) of the BROE server guarantees that the total execution needed by  $A_k$  to release the global lock will be granted with a delay of at most  $(P_k - \alpha_k P_k) = \Delta_k/2$  time-units. Therefore, the execution supplied to  $A_k$  in  $[t_l, t_{\text{miss}}]$  is at least  $(t_{\text{miss}} - t_l - \Delta_k/2)$ . Note that the execution supplied to an application does not depend on the particular local scheduling algorithm, but only on the global scheduling policy. Therefore, for both SRP and nonpreemptive locking,  $s(t_l, t_{\text{miss}}) \geq t_{\text{miss}} - t_l - \Delta_k/2$ . Equation (32) then becomes

$$W(Y) + \xi^{\text{NPL}}(t_s, t_{\text{miss}}) > t_{\text{miss}} - t_l - \frac{\Delta_k}{2} - s(t_l, t_s).$$

<sup>5</sup>We believe the resource holding time to be similar to the *maximum critical section execution time* defined in [40].

Using  $h_k(R_\ell) \geq \xi^{\text{NPL}}(t_l, t_s) + \xi^{\text{NPL}}(t_s, t_{\text{miss}})$ , it follows:

$$W(Y) + h_k(R_\ell) - \xi^{\text{NPL}}(t_l, t_s) > t_{\text{miss}} - t_l - \frac{\Delta_k}{2} - s(t_l, t_s).$$

Since the only job of  $A_k$  that is scheduled in  $[t_l, t_s]$  is the blocking job,  $\xi^{\text{NPL}}(t_l, t_s) \equiv s(t_l, t_s)$ . Then

$$W(Y) + h_k(R_\ell) > t_{\text{miss}} - t_l - \frac{\Delta_k}{2}.$$

Being  $h_k(R_\ell) \leq \Delta_k/2$ , the total amount of execution requested in  $[t_s, t_{\text{miss}}]$  by jobs in  $Y$  is

$$W(Y) > t_{\text{miss}} - t_l - \Delta_k.$$

Since we assumed that when  $A_k$  is scheduled on a dedicated processor of speed  $\alpha_k \leq 1$ , each job completes at least  $\Delta_k$  time-units before its deadline, the total demand of jobs in  $Y$  cannot be larger than the RHS of the above equation, reaching a contradiction and proving the theorem. ■

We believe Theorem 7 represents a very interesting result, since it allows improving the schedulability of the system both locally and globally. Notice, if the above theorem is satisfied for some  $A_k$  and  $R_\ell$ , then we may use  $h_k(R_\ell)$  instead of  $H_k(R_\ell)$  in the admission control tests of Section III-E. This increases the likelihood of  $A_k$  being admitted because the amount  $A_k$  could block applications  $A_i$  with  $P_i < P_k$  is decreased. Moreover, it allows using a very simple nonpreemptive locking protocol, decreasing the number of preemptions experienced by an application while holding a lock, and avoiding the use of more complex protocols for the access to shared resources.

## VII. SHARING GLOBAL RESOURCES

One of the features of our open environment that distinguishes it from other work that also considers resource-sharing is our approach towards the sharing of global resources across applications.

As stated above, there are works [18] mandating that global resources be accessed with local preemptions disabled. The rationale behind this approach is sound: by holding global resources for the least possible amount of time, each application minimizes the blocking interference to which it subjects other applications. However, the downside of such nonpreemptive execution is felt *within* each application—by requiring certain critical sections to execute nonpreemptively, it is more likely that an application when evaluated in isolation upon its slower-speed VP will be deemed infeasible. The server framework and analysis described in this paper allows for several possible execution modes for critical sections. We now analyze when each mode may be used.

More specifically, in extracting the interface for an application  $A_k$  that uses global resources, we can distinguish between three different cases.

- If the application is feasible on its VP when it executes a global resource  $R_\ell$  nonpreemptively, then have it execute  $R_\ell$  nonpreemptively.
- If an application is infeasible on its VP of speed  $\alpha_k$  when scheduled using EDF+SRP for  $R_\ell$ , it follows from the optimality of EDF+SRP [6] that no (work-conserving) scheduling strategy can result in this application being feasible upon a VP of the specified speed. Thus, by Theorem 5, no application-level scheduler can guarantee deadlines will be

met for the application on any BROE server with parameter  $\alpha_k$ .

- The interesting case is when neither of the two above holds: the system is infeasible when  $R_\ell$  executes nonpreemptively but feasible when access to  $R_\ell$  is arbitrated using the SRP. In that case, the objective should be to *devise a local scheduling algorithm for the application that retains feasibility while minimizing the resource holding times*. There are two possibilities.
  - a) Let  $h_k(R_\ell)$  be the largest critical section of any job of  $A_k$  that accesses global resource  $R_\ell$ . If  $h_k(R_\ell) \leq \Delta_k/2$  (in addition to the previously stated constraint on the resource-hold time  $H_k(R_\ell) \leq \alpha_k P_k$ ), then  $A_k$  may disable (local) preemptions when executing global resource  $R_\ell$  on its BROE server. In some cases, it may be advantageous to reduce  $H_k(R_\ell)$  to increase the chances that the constraint  $H_k(R_\ell) \leq \alpha_k P_k$  is satisfied.
  - b) If  $h_k(R_\ell) > \Delta_k/2$  but  $H_k(R_\ell) \leq \alpha_k P_k$  still holds,  $R_\ell$  may be executed using SRP. The resource-hold time could potentially be reduced by using techniques from [13], as discussed in Section VI.

## VIII. RELATION TO EXISTING WORKS

We consider this paper to be a generalization of earlier (“first-generation”) open environments (see, e.g., [11], [17], [21], [22], [35], and [41]), in that our results are applicable to shared platforms composed of serially reusable shared resources in addition to a preemptive processor. These projects assume that each individual application is composed of periodic implicit-deadline (“Liu and Layland”) tasks that do not share resources (neither locally within each application nor globally across applications); however, the resource “supply” models considered turn out to be alternative implementations of our scheduler (in the absence of shared resources).

There are other works that instead consider the access to globally shared resources in simpler server-based environments. The servers used in these environments do not allow for the hierarchical execution of different applications, but are often limited to the execution of one single task per server. Among fixed priority scheduled servers, Kuo and Li presented in [25] a resource sharing approach to be used with sporadic servers [43]. All global resources are handled by a single dedicated server that has capacity equal to the sum of all critical sections lengths and period equal to the GCD of the periods of the tasks accessing global resources. The drawback of this approach is that it could require a large utilization to accommodate such a dedicated server.

Among EDF-scheduled task systems, Ghazalie and Baker present in [24] an overrun mechanism to be used together with dynamic deferrable, sporadic and deadline exchange servers. Caccamo and Sha propose in [16] a modification to the CBS server [1] with a rule that allows locking a global critical section only when the CBS server has enough capacity to serve the whole critical section. Otherwise, the capacity is recharged and the server deadline is postponed. Since the deadline could be arbitrarily postponed by a task monopolizing the CPU (*deadline aging problem*), other tasks sharing the same server could potentially need to wait for an arbitrarily large amount of

time. Having “operative” deadline greater than the original ones can therefore violate the bounded-delay property, rendering the algorithm unsuitable for hierarchical open environments.

The Bandwidth Inheritance protocol (BWI) presented by Lipari *et al.* in [30] is based on a Priority Inheritance approach. A server holding a lock on a shared resource inherits the bandwidth of each server blocked on the same resource. The advantage of this protocol is that it does not require to know in advance which shared resources is accessed by each task, since it does not make use of the concept of “ceiling.” However, knowing the worst-case critical section lengths is necessary to compute the blocking due to shared resources for admission control. The computation of the interference due to other components is rather complex. Moreover, particular strategies should be used for deadlock avoidance when designing an application.

Feng presents in [21] two mechanisms for the access to shared resources in a server-based environment. In the first one, there is one server for each globally shared resource; whenever a job of an application requires access to a global resource, the job is executed [in first-in–first-out (FIFO) order] in the server associated with that resource, while the application can continue executing other jobs in the initial server. In the second approach, called “partition coalition,” a further mechanism is added: each task that may access a global resource is assigned a share of bandwidth (subpartition); every time a task blocks other tasks on a global resource, it inherits the bandwidth assigned to the blocked tasks, together with the (potentially null) bandwidth of the server associated to the global resource. Global critical sections are therefore executed in a “coalition” of partitions. Both approaches have some analogies with the “multireserve PCP” presented by de Niz *et al.* in [19], where instead of serving the blocked tasks in FIFO order, a priority-based approach with resource ceilings is used to limit the blocking times. Also, these works seem more suited for servers with one single task than for a hierarchical environment. Moreover, the schedulability analysis of such systems seems rather complex if compared with the admission control algorithms presented in Section III-E.

Among works that are closer in scope and ambition to this work, there is the solution proposed by Davis and Burns in [18], the work developed in Mälardalen [9], [10], [40], and the FIRST Scheduling Framework (FSF) [2]. Like these projects, our approach models each individual application as a collection of sporadic tasks which may share resources. One major difference between our work and most of these related works concerns the case in which the budget is exhausted while an application is still holding a global lock. The works in [2], [10], [18], and [40] introduced an overrun mechanism that continues executing the application until the lock is released, even if the budget is exhausted. The drawback of this approach is that in the scheduling analysis it is necessary to account for the largest overrun of each task/server, significantly reducing the available schedulable bandwidth.<sup>6</sup> We instead decided to start executing a global critical section only if the remained budget is sufficiently large to execute the whole critical section. Otherwise, we delay

<sup>6</sup>To reduce the overrun, Davis and Burns propose in [18] to execute each global critical section with local preemptions disabled. However, this imposes locally a larger interference on high priority jobs that do not access any global resource.

the acquisition of the global lock until the budget can be safely recharged.<sup>7</sup> We showed that in our framework this delay does not cause schedulability penalties.

Another difference between our work and the results presented in [18] concerns modularity. We have adopted an approach wherein each application is evaluated in isolation, and integration of the applications into the open environment is done based upon only the (relatively simple) interfaces of the applications. By contrast, [18] presents a monolithic approach to the entire system, with top-level schedulability formulas that cite parameters of individual tasks from different applications. We expect that a monolithic approach is more accurate but does not scale, and is not really in keeping with the spirit of open environment design.

As a final remark, it is possible to integrate our server with some previously proposed reclaiming mechanism to exploit the unused bandwidth. A simple rule can be easily added by setting to *Inactive* the state of all servers when the processor is idle. Moreover, the reclaiming mechanism used by GRUB in [27] may be implemented by updating the virtual time  $V_k$  of an executing application  $A_k$  at a rate  $\alpha_{\text{active}}/\alpha_k$ , instead than at a rate  $1/\alpha_k$ , where  $\alpha_{\text{active}}$  represents the sum of the  $\alpha_k$  of each admitted application  $A_k$  that is either in *Contending*, *Noncontending*, or *Suspended* state, i.e., excluding all inactive applications. We believe that other mechanisms, like the ones used in [14] and [15], can be adapted to the presented framework. We leave these problems for future works.

## IX. CONCLUSION AND FUTURE WORK

In this paper, we have presented a design for an open environment that allows for multiple independently developed and validated applications to be multiprogrammed on to a single shared platform. We believe that our design contains many significant innovations.

- We have defined a clean interface between applications and the environment, which encapsulates the important information while abstracting away unimportant details.
- The simplicity of the interface allows for efficient runtime admission control, and helps avoid combinatorial explosion as the number of applications increases.
- We have addressed the issue of interapplication resource sharing in great detail. moving beyond the ad hoc strategy of always executing shared global resources nonpreemptively, we have instead formalized the desired property of such resource-sharing strategies as *minimizing resource holding times*.
- We have studied a variety of strategies for performing arbitration for access to shared global resources within individual applications such that resource holding times are indeed minimized.

For the sake of concreteness, we have analyzed the local schedulability of individual applications that are executed upon our open environment using EDF and some protocol for arbitrating access to shared resources. This is somewhat constraining—ideally, we would like to be able to have each application scheduled using *any* local scheduling algorithm.<sup>8</sup>

<sup>7</sup>A similar design choice has been taken for the SIRAP protocol described in [9].

<sup>8</sup>Shin and Lee [41] refer to this property as *universality*.

Let us first address the issue of the task models that may be used in our approach. The results obtained in this paper have assumed that each application is composed of a collection of jobs that share resources. Therefore, the results contained in the paper extend in a straightforward manner to the situation where individual applications are represented using more general task models such as the multiframe [33], [34], generalized multiframe [7], or recurring [4], [5] task models—in essence, any formal model satisfying the *task independence assumptions* [7] may be used.

We conjecture that our framework can also handle applications modeled using task models not satisfying the task independence assumptions, provided the resource sharing mechanism used is independent of the absolute deadlines of the jobs, and only depends upon the relative priorities of the jobs according to EDF. We believe that our approach is general enough to successfully schedule such applications that have been validated by hand on a slower processor; we are currently working on proving this conjecture.

Next, let us consider local scheduling algorithms. We expect that analysis similar to ours could be conducted if a local application were to instead use (say) the deadline-monotonic scheduling algorithm [26], [31] with sporadic tasks, or some other fixed priority assignment with some more general task model (again, satisfying the task independence assumption). As discussed in Section IV-C, prior work on scheduling on resource partitions has assumed the local tasks do not share resources; we believe these results could be easily extended to include local resource sharing and used within our server framework.

A final note concerning generalizations. Our approach may also be applied to applications which are scheduled using *table-driven scheduling*, in which the entire sequence of jobs to be executed is precomputed and stored in a lookup table prior to runtime. Local scheduling for such systems reduces to dispatch based on table-lookup: such applications are also successfully scheduled by our open environment.

#### ACKNOWLEDGMENT

We are extremely grateful to the anonymous reviewers of both the conference and journal versions of this paper for their help in improving the quality of the paper.

#### REFERENCES

- [1] L. Abeni and G. Buttazzo, "Integrating multimedia applications in hard real-time systems," in *Proc. Real-Time Syst. Symp.*, Madrid, Spain, Dec. 1998, pp. 3–13.
- [2] M. Aldea, G. Bernat, I. Broster, A. Burns, R. Dobrin, J. M. Drake, G. Fohler, P. Gai, M. G. Harbour, G. Guidi, J. Gutiérrez, T. Lennvall, G. Lipari, J. Martínez, J. Medina, J. Palencia, and M. Trimarchi, "FSF: A real-time scheduling architecture framework," in *Proc. IEEE Real-Time Technol. Appl. Symp. (RTAS)*, Los Alamitos, CA, 2006, pp. 113–124.
- [3] T. P. Baker, "Stack-based scheduling of real-time processes," *Real-Time Systems: Int. J. Time-Critical Comput.*, vol. 3, pp. 67–99, 1991.
- [4] S. Baruah, "A general model for recurring real-time tasks," in *Proc. Real-Time Syst. Symp.*, Madrid, Spain, Dec. 1998, pp. 114–122.
- [5] S. Baruah, "Dynamic- and static-priority scheduling of recurring real-time tasks," *Real-Time Systems: Int. J. Time-Critical Comput.*, vol. 24, no. 1, pp. 99–128, 2003.
- [6] S. Baruah, "Resource sharing in EDF-scheduled systems: A closer look," in *Proc. IEEE Real-Time Syst. Symp.*, Rio de Janeiro, Dec. 2006, pp. 379–387.
- [7] S. Baruah, D. Chen, S. Gorinsky, and A. Mok, "Generalized multiframe tasks," *Real-Time Systems: Int. J. Time-Critical Comput.*, vol. 17, no. 1, pp. 5–22, Jul. 1999.
- [8] S. Baruah, A. K. Mok, and L. E. Rosier, "Preemptively scheduling hard-real-time sporadic tasks on one processor," in *Proc. 11th Real-Time Syst. Symp.*, Orlando, FL, 1990, pp. 182–190.
- [9] M. Behnam, I. Shin, T. Nolte, and M. Nolin, "SIRAP: A synchronization protocol for hierarchical resource sharing in real-time open systems," in *Proc. 7th ACM and IEEE Int. Conf. Embedded Softw.: EM-SOFT'07*, Salzburg, Austria, 2007.
- [10] M. Behnam, I. Shin, T. Nolte, and M. Nolin, "Scheduling of semi-independent real-time components: Overrun methods and resource holding times," in *Proc. 13th IEEE Int. Conf. Emerging Technol. Factory Autom. (ETFA'08)*, Hamburg, Germany, Sep. 2008, pp. 575–582.
- [11] G. Bernat and A. Burns, "Multiple servers and capacity sharing for implementing flexible scheduling," *Real-Time Syst.*, vol. 22, pp. 49–75, 2002.
- [12] M. Bertogna, N. Fisher, and S. Baruah, "Static-priority scheduling and resource hold times," in *Proc. Int. Workshop Parallel Distribut. Real-Time Syst. (IPDPS)*, Long Beach, CA, Mar. 2007, pp. 1–8.
- [13] M. Bertogna, N. Fisher, and S. Baruah, "Resource holding times: Computation and optimization," *Real-Time Syst.*, vol. 41, no. 2, pp. 87–117, Feb. 2009.
- [14] M. Caccamo, G. Buttazzo, and L. Sha, "Capacity sharing for overrun control," in *Proc. 21th IEEE RTSS*, Orlando, FL, 2000, pp. 295–304.
- [15] M. Caccamo, G. Buttazzo, and D. Thomas, "Efficient reclaiming in reservation-based real-time systems with variable execution times," *IEEE Trans. Comput.*, vol. 54, no. 2, pp. 198–213, Feb. 2005.
- [16] M. Caccamo and L. Sha, "Aperiodic servers with resource constraints," in *Proc. IEEE Real-Time Syst. Symp.*, London, U.K., Dec. 2001, pp. 161–170.
- [17] R. I. Davis and A. Burns, "Hierarchical fixed priority pre-emptive scheduling," in *Proc. IEEE Real-Time Syst. Symp.*, Miami, FL, 2005, pp. 389–398.
- [18] R. I. Davis and A. Burns, "Resource sharing in hierarchical fixed priority pre-emptive systems," in *Proc. IEEE Real-Time Syst. Symp.*, Rio de Janeiro, Brazil, 2006, pp. 257–268.
- [19] D. de Niz, L. Abeni, S. Saewong, and R. Rajkumar, "Resource sharing in reservation-based systems," in *Proc. IEEE Real-Time Syst. Symp.*, London, U.K., Dec. 2001, pp. 171–180.
- [20] Z. Deng and J. Liu, "Scheduling real-time applications in an open environment," in *Proc. 18th Real-Time Syst. Symp.*, San Francisco, CA, Dec. 1997, pp. 308–319.
- [21] X. Feng, "Design of real-time virtual resource architecture for large-scale embedded systems," Ph.D. dissertation, Dept. Comput. Sci., Univ. Texas at Austin, Austin, TX, 2004.
- [22] X. A. Feng and A. Mok, "A model of hierarchical real-time virtual resources," in *Proc. IEEE Real-Time Syst. Symp.*, 2002, pp. 26–35.
- [23] N. Fisher, M. Bertogna, and S. Baruah, "Resource-locking durations in EDF-scheduled systems," in *Proc. 13th IEEE Real-Time and Embedded Technol. Appl. Symp. (RTAS)*, Bellevue, WA, 2007, pp. 91–100.
- [24] T. M. Ghazalie and T. P. Baker, "Aperiodic servers in a deadline scheduling environment," *Real-Time Systems: Int. J. Time-Critical Comput.*, vol. 9, pp. 31–67, 1995.
- [25] T.-W. Kuo and C.-H. Li, "A fixed priority driven open environment for real-time applications," in *Proc. IEEE Real-Time Syst. Symp.*, Madrid, Spain, Dec. 1999, p. 256.
- [26] J. Y.-T. Leung and J. Whitehead, "On the complexity of fixed-priority scheduling of periodic, real-time tasks," *Perform. Eval.*, vol. 2, pp. 237–250, 1982.
- [27] G. Lipari and S. Baruah, "Greedy reclaimation of unused bandwidth in constant-bandwidth servers," in *Proc. EuroMicro Conf. Real-Time Syst.*, Stockholm, Sweden, Jun. 2000, pp. 193–200.
- [28] G. Lipari and E. Bini, "Resource partitioning among real-time applications," in *Proc. EuroMicro Conf. Real-Time Syst.*, Porto, Portugal, 2003, pp. 151–160.
- [29] G. Lipari and G. Buttazzo, "Schedulability analysis of periodic and aperiodic tasks with resource constraints," *J. Syst. Arch.*, vol. 46, no. 4, pp. 327–338, 2000.
- [30] G. Lipari, G. Lamastra, and L. Abeni, "Task synchronization in reservation-based real-time systems," *IEEE Trans. Comput.*, vol. 53, no. 12, pp. 1591–1601, 2004.
- [31] C. L. Liu and J. Layland, "Scheduling algorithms for multiprogramming in a hard real-time environment," *J. ACM*, vol. 20, no. 1, pp. 46–61, 1973.
- [32] A. K. Mok, "Fundamental design problems of distributed systems for the hard-real-time environment," Ph.D. dissertation, Lab. Comput. Sci., Massachusetts Institute of Technology, Cambridge, MA, 1983, Available as Tech. Rep. MIT/LCS/TR-297.
- [33] A. K. Mok and D. Chen, "A multiframe model for real-time tasks," in *Proc. 17th Real-Time Syst. Symp.*, Washington, DC, 1996, p. 22.

- [34] A. K. Mok and D. Chen, "A multiframe model for real-time tasks," *IEEE Trans. Softw. Eng.*, vol. 23, no. 10, pp. 635–645, Oct. 1997.
- [35] A. K. Mok, X. Feng, and D. Chen, "Resource partition for real-time systems," in *Proc. 7th IEEE Real-Time Technol. Appl. Symp. (RTAS'01)*, May 2001, pp. 75–84, IEEE.
- [36] R. Pellizzoni and G. Lipari, "Feasibility analysis of real-time periodic tasks with offsets," *Real-Time Systems: Int. J. Time-Critical Comput.*, vol. 30, no. 1–2, pp. 105–128, May 2005.
- [37] R. Rajkumar, *Synchronization in Real-Time Systems—A Priority Inheritance Approach*. Boston, MA: Kluwer, 1991.
- [38] R. Rajkumar, K. Juvva, A. Molano, and S. Oikawa, "Resource kernels: A resource-centric approach to real-time and multimedia systems," *Readings in Multimedia Computing and Networking*, pp. 476–490, 2001.
- [39] L. Sha, R. Rajkumar, and J. P. Lehoczky, "Priority inheritance protocols: An approach to real-time synchronization," *IEEE Trans. Comput.*, vol. 39, no. 9, pp. 1175–1185, 1990.
- [40] I. Shin, M. Behnam, T. Nolte, and M. Nolin, "Synthesis of optimal interfaces for hierarchical scheduling with resources," in *Proc. 29th IEEE Int. Real-Time Syst. Symp. (RTSS'08)*, Barcelona, Spain, Dec. 2008, pp. 209–220.
- [41] I. Shin and I. Lee, "Periodic resource model for compositional real-time guarantees," in *Proc. IEEE Real-Time Syst. Symp.*, 2003, pp. 2–13.
- [42] I. Shin and I. Lee, "Compositional real-time scheduling framework," in *Proc. IEEE Real-Time Syst. Symp.*, 2004, pp. 57–67.
- [43] B. Sprunt, L. Sha, and J. P. Lehoczky, "Scheduling sporadic and aperiodic events in a hard real-time system," *Real-Time Systems: Int. J. Time-Critical Comput.*, vol. 1, no. 1, pp. –, Jun. 1989, ESD-TR-89-19.



**Marko Bertogna** graduated (*summa cum laude*) in telecommunications engineering from the University of Bologna, Bologna, Italy, in 2002. He received the Ph.D. degree in computer science from Scuola Superiore Sant'Anna, Pisa, Italy, in 2008.

He is an Assistant Professor at the Scuola Superiore Sant'Anna. His research interests include scheduling and schedulability analysis of real-time multiprocessor systems, protocols for the exclusive access to shared resources, hierarchical systems, and reconfigurable devices.

Prof. Bertogna received the 2005 IEEE/Euromicro Conference on Real-Time Systems Best Paper Award.



**Nathan Fisher** received the B.S. degree from the University of Minnesota, Minneapolis, in 1999, the M.S. degree from Columbia University, New York, in 2002, and the Ph.D. degree from the University of North Carolina at Chapel Hill in 2007, all in computer science.

He is an Assistant Professor with the Department of Computer Science, Wayne State University. His research interests are in real-time and embedded computer systems, parallel and distributed algorithms, resource allocation, and approximation algorithms. His current research focus is on multiprocessor scheduling theory and composability of real-time applications.



**Sanjoy Baruah** received the Ph.D. degree from the University of Texas at Austin in 1993.

He is a Professor with the Department of Computer Science, University of North Carolina at Chapel Hill. His research and teaching interests are in scheduling theory, real-time and safety-critical system design, and resource-allocation and sharing in distributed computing environments.