

Supporting soft real-time parallel applications on multiprocessors[☆]Cong Liu^{a,*}, James H. Anderson^b^a Department of Computer Science, The University of Texas at Dallas, United States^b Department of Computer Science, The University of North Carolina at Chapel Hill, United States

ARTICLE INFO

Article history:

Available online 31 July 2013

Keywords:

Multiprocessor scheduling
Real-time systems
Parallel applications
Response time bounds

ABSTRACT

The prevalence of multicore processors has resulted in the wider applicability of parallel programming models such as OpenMP and MapReduce. A common goal of running parallel applications implemented under such models is to guarantee bounded response times while maximizing system utilization. Unfortunately, little previous work has been done that can provide such performance guarantees. In this paper, this problem is addressed by applying soft real-time scheduling analysis techniques. Analysis and conditions are presented for guaranteeing bounded response times for parallel applications under global EDF multiprocessor scheduling.

Published by Elsevier B.V.

1. Introduction

The growing prevalence of multicore platforms has resulted in the wider applicability of parallel programming models such as OpenMP [3] and MapReduce [5]. Such models can be applied to parallelize certain segments of programs, thus better utilizing hardware resources and possibly shortening response times. Many applications implemented under such parallel programming models have soft real-time (SRT) constraints. Examples include real-time parallel video and image processing applications [1,7] and computer vision applications such as colliding face detection and feature tracking [11]. In these applications, providing fast and bounded response times for individual video frames is important, to ensure smooth video output. However, achieving this at the expense of using conservative hard real-time (HRT) analysis is not warranted. In this paper, we consider how to schedule parallel task systems that require such SRT performance guarantees on multicore processors.

Parallel task models pose new challenges to real-time scheduling since intra-task parallelism has to be specifically considered. Recent papers [12,29] on scheduling real-time periodic parallel tasks have focused on providing HRT guarantees under global-earliest-deadline-first (GEDF) or partitioned deadline-monotonic (PDM) scheduling. However, as discussed above, viewing parallel tasks as HRT may be overkill in many settings and furthermore may result in significant schedulability-related utilization loss. Thus, our focus is to instead ensure bounded response times in

supporting parallel task systems by applying SRT scheduling analysis techniques. Specifically, we assign deadlines to parallel tasks and schedule them using GEDF, but in contrast to previous work [12,29], we allow deadlines to be missed provided the extent of such misses is bounded (hence response times are bounded as well). Moreover, we consider a generalized parallel task model that removes some of the restrictions seen in previous work (as discussed below).

Response time bounds have been studied extensively in the context of global real-time scheduling algorithms such as GEDF [6,13–25]. It has been shown that a variety of such algorithms can ensure bounded response times in ordinary real-time sporadic task systems (i.e., without intra-task parallelism) with no utilization loss on multiprocessors [6,13].¹ Motivated by these results, we consider whether it is possible to specify reasonable constraints under which bounded response times can be guaranteed using global real-time scheduling techniques, for sporadic parallel task systems that are not HRT in nature.

Related work. Scheduling non-real-time parallel applications is a deeply explored topic [4,5,8,9,26,31,32]. However, in most (if not all) prior work on this topic, including all of the just-cited work, scheduling decisions are made on a best-effort basis, so none of these results can provide performance guarantees such as response time bounds.

Regarding scheduling HRT parallel task systems, Lakshmanan et al. proposed a scheduling technique for the *fork-join* model, where a parallel task is a sequence of segments, alternating between sequential and parallel phases [12]. A sequential phase

[☆] Work supported by AT&T, IBM, Intel, and Sun Corps.; NSF Grants CNS 0834270, CNS 0834132, and CNS 0615197; and ARO Grant W911NF-06-1-0425.

* Corresponding author. Tel.: +1 9193601521; fax: +1 9199621799.

E-mail address: cong.liu@utdallas.edu (C. Liu).

¹ Technically, bounded response times can only be ensured for task systems that do not over-utilize the underlying platform. In all claims in this paper concerning bounded response times, a non-over-utilized system is assumed.

contains only one thread while a parallel phase contains multiple threads that can be executed concurrently on different processors. In their model, all parallel phases are assumed to have the same number of parallel threads, which must be no greater than the number of processors. Also, all threads in any parallel segment must have the same execution cost. The authors derived a resource augmentation bound of 3.42 under PDM scheduling.

In [29], Saifullah et al. extended the fork-join model so that each parallel phase can have a different number of threads and threads can have different execution costs. The authors proposed an approach that transforms each periodic parallel task into a number of ordinary constrained-deadline periodic tasks by creating per segment intermediate deadlines. They also showed that resource augmentation bounds of 2.62 and 3.42 can be achieved under GEDF and PDM scheduling, respectively. In [27], Nelissen et al. proposed techniques that optimize the number of processors needed to schedule sporadic parallel tasks. The authors also proved that the proposed techniques achieve a resource augmentation bound of 2.0 under scheduling algorithms such as U-EDF [28] and PD² [30].

In this paper, we seek to efficiently support parallel task systems on multiprocessors with bounded response times. We consider the general parallel task model as presented in [27,29]. A fundamental difference between this work and prior work is that we propose a SRT schedulability analysis framework to derive conditions for guaranteeing bounded response times.

Contributions. In this paper, we show that by assigning deadlines to parallel task systems and scheduling them under GEDF, such systems can be supported on multiprocessors with bounded response times. Our analysis shows that on a two-processor platform, no utilization loss results for any parallel task system. Despite this special case, on a platform with more than two processors, utilization constraints are needed. To discern how severe such constraints must fundamentally be, we present a parallel task set with minimum utilization that is unschedulable on any number of processors. This task set violates our derived constraint and has unbounded response times. The impact of utilization constraints can be lessened by restructuring tasks to reduce intra-task parallelism. We propose optimization techniques that can be applied to determine such a restructuring. Finally, we present the results of experiments conducted to evaluate the applicability of the derived schedulability condition.

Organization. The rest of this paper is organized as follows. Section 2 describes our system model. In Section 3, we present our analytical results. In Section 4, we discuss the above mentioned optimization technique. In Section 5, we experimentally evaluate the proposed analysis. Section 6 concludes.

2. System model

We consider the problem of scheduling a set $\tau = \{\tau_1, \dots, \tau_n\}$ of n independent sporadic parallel tasks on m processors. Each parallel task τ_i is a sequence of s_i segments, where the j th segment τ_i^j contains a set of v_i^j threads ($v_i^j > m$ is allowed). The k th ($1 \leq k \leq v_i^j$) thread $\tau_{i,k}^{j,k}$ in segment τ_i^j has a worst-case execution time of $e_{i,k}^{j,k}$. We assume that each thread $\tau_{i,k}^{j,k}$ executes for exactly $e_{i,k}^{j,k}$ time units. This assumption can be eased to treat $e_{i,k}^{j,k}$ as an upper bound, at the expense of more cumbersome notation. For notational convenience, we order the threads of each segment τ_i^j of each parallel task τ_i in largest-worst-case-execution-time-first order. Thus, thread $\tau_{i,1}^{j,1}$ has the largest worst-case execution time among all threads in any segment τ_i^j . For any segment τ_i^j , if $v_i^j > 1$, then the threads in this segment can be executed in parallel on different processors. The threads in the j th segment can execute only after all threads of $(j-1)$ th segment (if any) have completed. We let v_i^{\max} denote the maximum number of threads in any segment of task τ_i . We assume $v_i^{\max} \geq 2$ holds for at least one task

τ_i ; otherwise, the considered task system is simply an ordinary sporadic task system (without intra-task parallelism).

The worst-case execution time of any segment τ_i^j is defined as $e_i^j = \sum_{k=1}^{v_i^j} e_{i,k}^{j,k}$ (when all threads execute sequentially). The worst-case execution time of any parallel task τ_i is defined as $e_i = \sum_{j=1}^{s_i} e_i^j$ (when all threads in each segment of the task execute sequentially). In our analysis, we also make use of the best-case execution time of τ_i on m processors (when τ_i is the only task executing on m processors), denoted e_i^{\min} . In general, for any parallel task τ_i , if we allow $v_i^{\max} \geq m$ and threads in each segment have different execution costs, then the problem of calculating e_i^{\min} is equivalent to the problem of *minimum makespan scheduling* [10], where we treat each thread in a segment as an independent job and seek to obtain the minimum completion time for executing all such jobs on m processors. This gives us per segment best-case execution times, which can be summed to yield e_i^{\min} . Unfortunately, this problem has been proven to be NP-hard [10]. This problem can be solved using a classical dynamic programming-based algorithm [10], which has exponential time complexity with respect to the per segment thread count. However, for some special cases where certain restrictions on the task model apply, we can easily calculate e_i^{\min} in linear time. For example, when $v_i^{\max} \leq m$ holds, $e_i^{\min} = \sum_{j=1}^{s_i} e_{i,1}^{j,1}$ since in this case all threads of each segment of τ_i can be executed in parallel on m processors and thread $\tau_{i,1}^{j,1}$ has the largest execution cost in each segment τ_i^j . Moreover, when all threads in each segment have equal execution costs, $e_i^{\min} = \sum_{j=1}^{s_i} \sum_{k=1}^{\lceil v_i^j/m \rceil} e_{i,k}^{j,1}$, because the execution of each segment τ_i^j can be viewed as the executions of $\lceil v_i^j/m \rceil$ sequential sub-segments, each with an equal execution cost of $e_{i,1}^{j,1}$.

Each parallel task is released repeatedly, with each such invocation called a *job*. The k th job of τ_i , denoted $\tau_{i,k}$, is released at time $r_{i,k}$. Associated with each task τ_i is a period p_i , which specifies the minimum time between two consecutive job releases of τ_i . We require $e_i^{\min} \leq p_i$ for any task τ_i ; otherwise, response times (defined next) can grow unboundedly. The utilization of a task τ_i is defined as $u_i = e_i/p_i$, and the utilization of the task system τ as $U_{\text{sum}} = \sum_{\tau_i \in \tau} u_i$. We require $U_{\text{sum}} \leq m$; otherwise, response times can grow unboundedly. For any job $\tau_{i,k}$ of task τ_i , its u th segment is denoted $\tau_{i,k}^u$, and the v th thread of this segment is denoted $\tau_{i,k}^{u,v}$. An example parallel task is shown in Fig. 1. For clarity, a summary of important terms defined so far, as well as some additional terms defined later, is presented in Table 1.

Successive jobs of the same task are required to execute in sequence. If a job $\tau_{i,k}$ completes at time t , then its *response time* is $t - r_{i,k}$. A task's response time is the maximum response time of any of its jobs. Note that, when a job of a task completes after the release time of the next job of that task, this release time is not altered.

Assigning deadlines and priority points. Each parallel task τ_i has a specified relative deadline of d_i , which may differ from p_i (thus, our analysis is applicable to soft real-time arbitrary-deadline sporadic parallel tasks). We do not use such deadlines in prioritizing jobs, but rather assign each job $\tau_{i,k}$ a priority point at $d_{i,k} = r_{i,k} + p_i$ and schedule jobs on a global earliest-priority-point-first (GEPPF) basis. That is, earlier priority points are prioritized over later ones.² We assume that ties are broken by task ID (lower IDs are favored).

3. Response time bound

We derive a response time bound for GEPPF by comparing the allocations to a task system τ in a processor sharing (PS) schedule and an actual GEPPF schedule of interest for τ , both on m

² GEDF becomes a special case of GEPPF when $d_i = p_i$ holds for each τ_i .

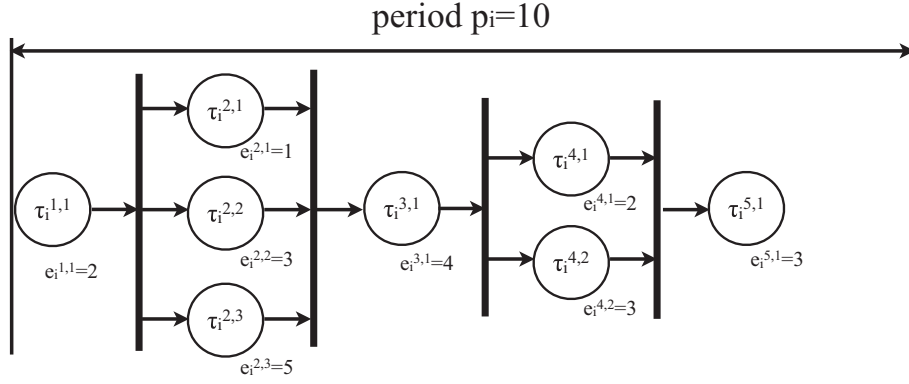


Fig. 1. Example parallel task τ_i . It has five segments where the second and fourth segments are parallel segments and contain three and two threads, respectively. This task has a worst-case execution cost of 23 time units, a period of 10 time units, and thus a utilization of 2.3.

Table 1
Summary of notation.

$\tau_{i,h}^j$	j th segment of the h th job of task τ_i
$\tau_{i,h}^{j,k}$	k th thread of segment $\tau_{i,h}^j$ of the h th job of task τ_i
s_i	Number of segments of task τ_i
$e_i^{j,k}$	Worst-case execution cost of thread $\tau_{i,h}^{j,k}$
e_i^j	Worst-case execution cost of segment $\tau_{i,h}^j$
e_i	Worst-case execution cost of task τ_i
e_i^{\min}	Best-case execution cost of task τ_i
ν_i^{\max}	Maximum number of threads in any segment of task τ_i
ν_{\max}	Maximum number of threads of any segment of the task that has the i th maximum number of threads of any segment among all tasks

processors, and quantifying the difference between the two. For any given sporadic parallel task system, a PS schedule is an ideal schedule where each released job $\tau_{i,k}$ executes with a rate equal to u_i (which ensures that each job completes exactly at its priority point). Note that parallelism is not considered in the PS schedule. A valid PS schedule exists for τ if $U_{\text{sum}} \leq m$ holds.

We analyze task allocations on a per-task basis.³

We assume time is discrete. For any time $t > 0$, the notation t^- is used to denote the time $t - \varepsilon$ in the limit $\varepsilon \rightarrow 0+$, and the notation t^+ is used to denote the time $t + \varepsilon$ in the limit $\varepsilon \rightarrow 0+$.

Definition 1. A task τ_i is *active* at time t if there exists a job $\tau_{i,h}$ such that $r_{i,h} \leq t < d_{i,h}$.

Definition 2. Job $\tau_{i,h}$ is *pending* at time t if $t \geq r_{i,h}$ and $\tau_{i,h}$ has not completed by t .

Definition 3. Job $\tau_{i,h}$ is *enabled* at t if $t \geq r_{i,h}$, $\tau_{i,h}$ has not completed by t , and $\tau_{i,h-1}$ (if $h > 1$) has completed by t . Similarly, any thread in segment $\tau_{i,h}^k$ is *enabled* at t if $t \geq r_{i,h}$, the thread has not completed by t , and all threads in segment $\tau_{i,h-1}^{k-1}$ (if any) have completed by t .

Let $A(\tau_{i,j}, t_1, t_2, S)$ denote the total allocation to the job $\tau_{i,j}$ in an arbitrary schedule S in $[t_1, t_2]$. Then, the total time allocated to all jobs of τ_i in $[t_1, t_2]$ in S is given by

$$A(\tau_i, t_1, t_2, S) = \sum_{j \geq 1} A(\tau_{i,j}, t_1, t_2, S).$$

³ The SRT analysis framework used here has been adopted from a framework for ordinary sporadic task systems first proposed in [6], and subsequently used in several other papers [13,19,24].

Consider a PS schedule PS . In such a schedule, τ_i executes with the rate u_i when it is active. (Recall that intra-task parallelism is not considered in the PS schedule.) Thus, if τ_i is active throughout $[t_1, t_2]$, then

$$A(\tau_{i,j}, t_1, t_2, PS) = (t_2 - t_1)u_i. \quad (1)$$

Note that according to the parallel task model, the term u_i in (1) could be greater than one. This is a key difference in comparison to most prior work where a PS schedule is considered. A PS schedule for an example task system is shown in Fig. 2.

The difference between the allocation to a job $\tau_{i,j}$ up to time t in a PS schedule and an arbitrary schedule S , denoted the *lag of job* $\tau_{i,j}$ at time t in schedule S , is defined by $\text{lag}(\tau_{i,j}, t, S) = A(\tau_{i,j}, 0, t, PS) - A(\tau_{i,j}, 0, t, S)$. The *lag of a task* τ_i at time t in schedule S is given by

$$\text{lag}(\tau_i, t, S) = \sum_{j \geq 1} \text{lag}(\tau_{i,j}, t, S) = A(\tau_i, 0, t, PS) - A(\tau_i, 0, t, S). \quad (2)$$

The concept of lag is important because, if lags remain bounded, then response times are bounded as well. The *LAG* for a finite job set J at time t in the schedule S is defined as

$$\begin{aligned} \text{LAG}(J, t, S) &= \sum_{\tau_{i,j} \in J} \text{lag}(\tau_{i,j}, t, S) \\ &= \sum_{\tau_{i,j} \in J} (A(\tau_{i,j}, 0, t, PS) - A(\tau_{i,j}, 0, t, S)). \end{aligned} \quad (3)$$

Our response time bound derivation focuses on a given task system τ . We order jobs in τ by EDF, and break ties by task ID. Let $\tau_{i,j}$ be a job of a task τ_i in τ , $t_d = d_{i,j}$, and S be a GEPPF schedule for τ with the following property.

(P) The response time of every job $\tau_{i,k}$ of higher priority than $\tau_{i,j}$ is at most $x + p_i + e_i$ in S , where $x \geq 0$.

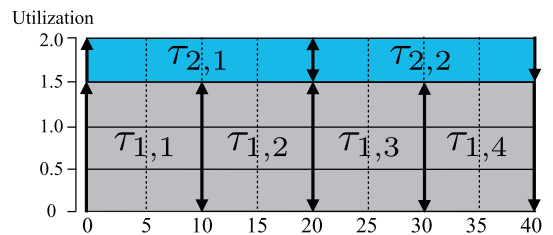


Fig. 2. PS schedule for a task system containing two tasks. Task τ_1 has a period of 10 time units and a utilization of 1.5. Task τ_2 has a period of 20 time units and a utilization of 0.5. As seen in the PS schedule, intra-task parallelism is not considered and each job completes exactly at its deadline.

Our objective is to determine the smallest x such that the response time of τ_{ij} is at most $x + p_l + e_l$. This would by induction imply a response time of at most $x + p_i + e_i$ for all jobs of every task τ_i , where $\tau_i \in \tau$. We assume that τ_{ij} finishes after t_d , for otherwise, its response time is trivially no greater than p_l . The steps for determining the value for x are as follows.

1. Determine an upper bound on the work pending for tasks in τ that can compete with τ_{ij} after t_d . This is dealt with in [Lemmas 1 and 2](#) in Section 3.1.
2. Determine a lower bound on the amount of work pending for tasks in τ that can compete with τ_{ij} after t_d , required for the response time of τ_{ij} to exceed $x + p_l + e_l$. This is dealt with in [Lemma 3](#) in Section 3.2.
3. Determine the smallest x such that the response time of τ_{ij} is at most $x + p_l + e_l$, using the above upper and lower bounds. This is dealt with in [Theorem 1](#) in Section 3.3.

Definition 4. $\mathbf{d} = \{\tau_{i,h} : (d_{i,h} < t_d) \vee (d_{i,h} = t_d \wedge i \leq l)\}$.

\mathbf{d} is the set of jobs with deadlines at most t_d with priority at least that of τ_{ij} . These jobs do not execute beyond t_d in the PS schedule. Note that τ_{ij} is in \mathbf{d} . Also note that jobs not in \mathbf{d} have lower priority than those in \mathbf{d} and thus do not affect the scheduling of jobs in \mathbf{d} . For simplicity, we will henceforth assume that no job not in \mathbf{d} executes in either the PS or GEPPF schedule. To avoid distracting “boundary cases,” we also assume that the schedule being analyzed is prepended with a schedule in which no deadlines are missed that is long enough to ensure that all previously released jobs referenced in the proof exist.

According to Property (P), job τ_{ij-1} has a response time of at most $x + p_l + e_l$. Thus, the completion time of τ_{ij-1} , denoted t_p (p for predecessor), is given by

$$t_p \leq r_{ij-1} + p_l + x + e_l \leq r_{ij} + x + e_l = t_d - p_l + x + e_l. \quad (4)$$

Definition 5. A time instant t is *busy* for a job set J if all m processors execute jobs in J at t . A time interval is busy for J if each instant within it is busy for J .

The following claim follows from the definition of LAG.

Claim 1. If $LAG(\mathbf{d}, t_2, S) > LAG(\mathbf{d}, t_1, S)$, where $t_2 > t_1$, then $[t_1, t_2]$ is non-busy for \mathbf{d} . In other words, LAG for \mathbf{d} can increase only throughout a non-busy interval.

An interval could be non-busy for \mathbf{d} only if there are not enough enabled jobs in \mathbf{d} to occupy all available processors.

Since \mathbf{d} includes all jobs of higher priority than τ_{ij} , the competing work for τ_{ij} after time t_d is given by the amount of work pending at t_d for jobs in \mathbf{d} , which is given by $LAG(\mathbf{d}, t_d, S)$.

3.1. Upper bound

In this section, we determine an upper bound on $LAG(\mathbf{d}, t_d, S)$. We first upper bound $lag(\tau_i, t, S)$ ($t \in [0, t_d]$) in [Lemma 1](#) below. Then, in [Lemma 2](#), we upper bound $LAG(\mathbf{d}, t_d, S)$ by summing individual task lags.

Definition 6. Let t_n be the end of the latest non-busy interval for \mathbf{d} before t_d , if any; otherwise, let $t_n = 0$ (see in [Fig. 3](#)).

By the above definition and [Claim 1](#), we have

$$LAG(\mathbf{d}, t_d, S) \leq LAG(\mathbf{d}, t_n, S). \quad (5)$$

Lemma 1. $lag(\tau_i, t, S) \leq u_i \cdot x + (u_i + 1) \cdot e_i$ for any task τ_i and $t \in [0, t_d]$.

Proof. Let $d_{i,k}$ be the deadline of the earliest pending job of τ_i , $\tau_{i,k}$, in the schedule S at time t . If such a job does not exist, then $lag(\tau_i, t, S) = 0$, and the lemma holds trivially. Let γ_i be the amount of work $\tau_{i,k}$ performs before t .

By the selection of $\tau_{i,k}$, we have

$$\begin{aligned} lag(\tau_i, t, S) &= \sum_{h \geq k} lag(\tau_{i,h}, t, S) \\ &= A(\tau_{i,k}, r_{i,k}, t, PS) - A(\tau_{i,k}, r_{i,k}, t, S) \\ &\quad + \sum_{h > k} (A(\tau_{i,h}, r_{i,h}, t, PS) - A(\tau_{i,h}, r_{i,h}, t, S)). \end{aligned} \quad (6)$$

By the definition of PS , $A(\tau_{i,k}, r_{i,k}, t, PS) \leq e_i$, and $\sum_{h > k} A(\tau_{i,h}, r_{i,h}, t, PS) \leq u_i \cdot \max(0, t - d_{i,k})$ (the latter follows because each such job $\tau_{i,h}$ executes with rate u_i in PS while active, and the sum of the active intervals under consideration is at most $t - d_{i,k}$). By the selection of $\tau_{i,k}$, $A(\tau_{i,k}, r_{i,k}, t, S) = \gamma_i$, and $\sum_{h > k} A(\tau_{i,h}, r_{i,h}, t, S) = 0$. By setting these values into (6), we have

$$lag(\tau_i, t, S) \leq e_i - \gamma_i + u_i \cdot \max(0, t - d_{i,k}). \quad (7)$$

There are two cases to consider.

Case 1. $d_{i,k} \geq t$. In this case, (7) implies $lag(\tau_i, t, S) \leq e_i - \gamma_i \leq u_i \cdot x + (u_i + 1) \cdot e_i$.

Case 2. $d_{i,k} < t$. In this case, because $t \leq t_d$ and $d_{i,j} = t_d$, $\tau_{i,k}$ is not the job τ_{ij} . Thus, by Property (P), $\tau_{i,k}$ has a response time of at most $x + p_i + e_i$. Since $\tau_{i,k}$ is the earliest pending job of τ_i at time t , the earliest possible completion time of $\tau_{i,k}$ is at t^+ . Thus, we have $t - r_{i,k} < t^+ - r_{i,k} \leq x + p_i + e_i$, which (because $d_{i,k} = r_{i,k} + p_i$) implies $t - d_{i,k} = t - r_{i,k} - p_i < x + e_i$.

Setting this value into (7), we have $lag(\tau_i, t, S) < e_i - \gamma_i + u_i \cdot (x + e_i) \leq u_i \cdot x + (u_i + 1) \cdot e_i$. \square

[Lemma 2](#) below upper bounds $LAG(\mathbf{d}, t_d, S)$. We first define some needed terms.

Definition 7. Let U be the sum of the $\min(m - 1, n)$ largest task utilizations. Let E be the largest value of the expression $\sum_{\tau_i \in \gamma} ((u_i + 1) \cdot e_i)$, where γ denotes any set of $\min(m - 1, n)$ tasks in τ .

Lemma 2. $LAG(\mathbf{d}, t_d, S) \leq U \cdot x + E$.

Proof. By (5), we have $LAG(\mathbf{d}, t_d, S) \leq LAG(\mathbf{d}, t_n, S)$. By summing individual task lags at t_n , we can bound $LAG(\mathbf{d}, t_n, S)$. If $t_n = 0$, then $LAG(\mathbf{d}, t_n, S) = 0$, so assume $t_n > 0$. Consider the set of tasks $\beta = \{\tau_i : \exists \tau_{i,h} \text{ such that } \tau_{i,h} \text{ is enabled at } t_n^-\}$. Given that the instant t_n^- is non-busy, there are not enough enabled jobs in \mathbf{d} to occupy all m processors. More precisely, there are not enough enabled threads belonging to jobs in \mathbf{d} to occupy all m processors. There could be at most $\min(m - 1, n)$ parallel tasks that have enabled jobs at t_n^- since each such parallel task has at least one enabled thread at t_n^- ; that is, $|\beta| \leq \min(m - 1, n)$.

If task τ_i does not have pending jobs at t_n^- , then $lag(\tau_i, t_n, S) \leq 0$. Therefore, we have

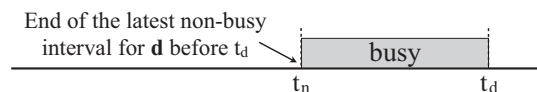


Fig. 3. Definition of t_n .

$$\begin{aligned}
& \text{LAG}(\mathbf{d}, t_d, S) \\
& \{ \text{by}(5) \} \\
& \leq \text{LAG}(\mathbf{d}, t_n, S) \\
& \{ \text{by}(3) \} \\
& = \sum_{\tau_i: \tau_{i,h}^w \in \mathbf{d}} \text{lag}(\tau_i, t_n, S) \\
& \leq \sum_{\tau_i \in \beta} \text{lag}(\tau_i, t_n, S) \\
& \{ \text{by Lemma 1} \} \\
& \leq \sum_{\tau_i \in \beta} (u_i \cdot x + (u_i + 1) \cdot e_i).
\end{aligned}$$

By Definition 7 and because $|\beta| \leq \min(m-1, n)$, we have $\text{LAG}(\mathbf{d}, t_d, S) \leq \sum_{\tau_i \in \beta} (u_i \cdot x + (u_i + 1) \cdot e_i) \leq U \cdot x + E$. \square

3.2. Lower bound

In the following lemma, we determine a lower bound on $\text{LAG}(\mathbf{d}, t_d, S)$ that is necessary for the response time of $\tau_{l,j}$ to exceed $x + p_l + e_l$.

Definition 8. If any thread of any segment of job $\tau_{i,h}$ is enabled at time t but does not execute at t , and at least one processor is executing some job other than $\tau_{i,h}$ at t , then $\tau_{i,h}$ is *preempted* at t (see Fig. 4).

Definition 9. Let v_{\max_i} denote the maximum number of threads of any segment of the task that has the i th maximum number of threads of any segment among tasks in τ .

If $\sum_{i=1}^n v_{\max_i} \leq m$, then each thread of each segment of each task in τ can be executed on a processor without being preempted, which implies that each task $\tau_k \in \tau$ has a bounded response time of $e_k^{\min} < x + p_k + e_k$. Thus, we consider the other case, where $\sum_{i=1}^n v_{\max_i} > m$. Moreover, since we assume that there exists at least one task $\tau_k \in \tau$ with $v_k^{\max} \geq 2$ (as discussed in Section 2), we have $v_{\max_1} \geq 2$. Thus, if $n > m$, then $\sum_{i=1}^m v_{\max_i} > m$ holds. Therefore, we have

$$\sum_{i=1}^{\min(m,n)} v_{\max_i} > m. \quad (8)$$

Definition 10. Let

$$Q = \begin{cases} 2 & \text{if } v_{\max_1} > m, \\ \min \left\{ k \mid \sum_{i=1}^k v_{\max_i} > m \right\} & \text{if } v_{\max_1} \leq m. \end{cases}$$

Q is used in Lemma 3 below to obtain a lower bound on $\text{LAG}(\mathbf{d}, t_d, S)$; the two conditions in the definition of Q arise because of different

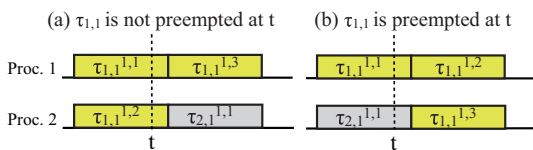


Fig. 4. Illustration of a preemption. Job $\tau_{1,1}$ has one segment with three parallel threads, executed on two processors. In inset (a), although $\tau_{1,1}^{1,3}$ is enabled but does not execute at time t , $\tau_{1,1}$ is not preempted at t since both processors are executing threads of $\tau_{1,1}$. In inset (b), $\tau_{1,1}$ is preempted by $\tau_{2,1}$ at t .

subcases considered in the proof of Lemma 3. Note that by the above definition and (8), we have

$$2 \leq Q \leq \min(m, n) \leq m. \quad (9)$$

Lemma 3. If the response time of $\tau_{l,j}$ exceeds $x + p_l + e_l$, then $\text{LAG}(\mathbf{d}, t_d, S) > Q \cdot x - (m-1) \cdot e_l$.

Proof. Throughout the proof of this lemma, we assume $\sum_{i=1}^n v_{\max_i} > m$ and $v_{\max_1} \geq 2$ both hold, for reasons discussed above. We prove the contrapositive: we assume that

$$\text{LAG}(\mathbf{d}, t_d, S) \leq Q \cdot x - (m-1) \cdot e_l \quad (10)$$

holds and show that the response time of $\tau_{l,j}$ cannot exceed $x + p_l + e_l$. Let η_l be the amount of work $\tau_{l,j}$ performs by time t_d in S . Define y as follows.

$$y = \frac{Q}{m} \cdot x + \frac{\eta_l}{m} \quad (11)$$

Let W be the amount of work due to jobs in \mathbf{d} that can compete with $\tau_{l,j}$ after $t_d + y$, including the work due for $\tau_{l,j}$. Let t_f be the completion time of $\tau_{l,j}$. We consider two cases.

Case 1. $[t_d, t_d + y)$ is a busy interval for \mathbf{d} . In this case, we have

$$\begin{aligned}
W &= \text{LAG}(\mathbf{d}, t_d, S) - my \\
& \{ \text{by}(10) \} \\
& \leq Q \cdot x - (m-1) \cdot e_l - my \\
& \{ \text{by}(11) \} \\
& = Q \cdot x - (m-1) \cdot e_l - Q \cdot x - \eta_l \\
& = -(m-1) \cdot e_l - \eta_l \\
& < 0.
\end{aligned}$$

Because GEPPF is work-conserving (i.e., GEPPF idles a processor only when there is no enabled job), at least one processor is busy until $\tau_{l,j}$ completes. Thus, the amount of work performed by the system for jobs in \mathbf{d} during the interval $[t_d + y, t_f]$ is at least $t_f - t_d - y$. Hence, $t_f - t_d - y \leq W < 0$. Therefore, the response time of $\tau_{l,j}$ is

$$\begin{aligned}
t_f - r_{l,j} &= t_f - t_d + p_l \\
& < y + p_l \\
& \{ \text{by}(11) \} \\
& = \frac{Q}{m} \cdot x + \frac{\eta_l}{m} + p_l \\
& \{ \text{by}(9) \} \\
& \leq x + e_l + p_l.
\end{aligned}$$

Case 2. $[t_d, t_d + y)$ is a non-busy interval for \mathbf{d} . Let $t_s \geq t_d$ be the earliest non-busy instant in $[t_d, t_d + y)$. Recall (see (4)) that t_p is the completion time of job $\tau_{l,j-1}$. We consider three subcases.

Subcase 2.1. $t_p \leq t_s$ and $\tau_{l,j}$ is not preempted within $[t_p, t_s)$. As illustrated in Fig. 5, in this case, $\tau_{l,j}$ can start execution at t_s because t_s is non-busy. Since $\tau_{l,j}$ is not preempted within $[t_s, t_p)$, $\tau_{l,j}$ completes by $t_s + e_l - \eta_l$. Thus, because $t_s < t_d + y$, $\tau_{l,j}$ finishes by time

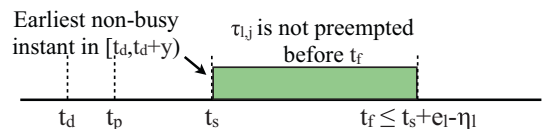


Fig. 5. Subcase 2.1.

$$\begin{aligned}
t_s + e_l - \eta_l &< t_d + y + e_l - \eta_l \\
&\{ \text{by(11)} \} \\
&= t_d + \frac{Q}{m} \cdot x + \frac{\eta_l}{m} + e_l - \eta_l \\
&\{ \text{by(9)} \} \\
&\leq r_{lj} + p_l + x + e_l.
\end{aligned}$$

Subcase 2.2. $t_p \leq t_s$ and τ_{lj} is preempted within $[t_p, t_s)$. If $t_f \leq y + t_d$, then

$$\begin{aligned}
t_f - r_{lj} &\leq y + t_d - r_{lj} \\
&\{ \text{by(11)} \} \\
&= \frac{Q}{m} \cdot x + \frac{\eta_l}{m} + p_l \\
&\{ \text{by(9)} \} \\
&\leq x + e_l + p_l.
\end{aligned}$$

So assume $t_f > y + t_d$. Let $t_1 > t_s$ be the earliest time when τ_{lj} is preempted. As shown in Fig. 6, by the definition of t_s and t_1 , τ_{lj} executes throughout $[t_s, t_1)$ without being preempted. Because τ_{lj} is preempted at t_1 , t_1 is busy with respect to \mathbf{d} . Let t_2 be the last time τ_{lj} resumes execution after being preempted if such a time exists; if such a time does not exist, which implies that τ_{lj} is preempted until t_f , then let $t_2 = t_f$ (note that by Definition 8, some threads of τ_l^j can execute while τ_l^j is preempted). Within $[t_1, t_2)$, τ_{lj} could be preempted multiple times. By Definition 8, all such intervals during which τ_{lj} is preempted must be busy in order for the preemption to happen. Given that $t_f \leq t_2 + e_l - \eta_l$, if $t_2 \leq y + t_d$, then $t_f \leq y + t_d + e_l - \eta_l$, in which case, because $t_d - r_{lj} = p_l$, the response time of τ_{lj} is

$$\begin{aligned}
t_f - r_{lj} &\leq y + p_l + e_l - \eta_l \\
&\{ \text{by(11)} \} \\
&\leq \frac{Q}{m} \cdot x + p_l + e_l \\
&\{ \text{by(9)} \} \\
&\leq x + p_l + e_l,
\end{aligned}$$

as required.

If $t_2 > t_d + y$, then the amount of work due to \mathbf{d} performed within $[t_d, t_d + y)$ is at least $my - (m - 1) \cdot \min(e_l, y)$ because all intervals during which τ_{lj} is preempted are busy, and τ_{lj} can execute for at most e_l time in $[t_d, y + t_d)$. (Within intervals in $[t_s, t_d + y)$ where τ_{lj} is not preempted, at least one processor is occupied by τ_{lj} .) Thus, the amount of work that can compete with τ_{lj} after $t_d + y$ is

$$\begin{aligned}
W &\leq \text{LAG}(\mathbf{d}, t_d, S) - (my - (m - 1) \cdot \min(e_l, y)) \\
&\{ \text{by(10)} \} \\
&\leq Q \cdot x - (m - 1) \cdot e_l - (my - (m - 1) \cdot \min(e_l, y)) \\
&\leq Q \cdot x - my \\
&\{ \text{by(11)} \} \\
&= -\eta_l \\
&\leq 0.
\end{aligned}$$

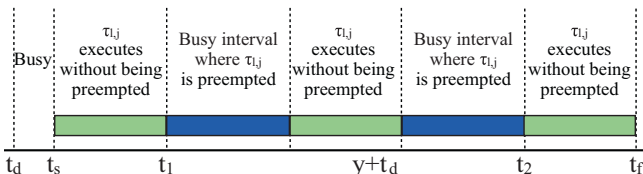


Fig. 6. Subcase 2.2.

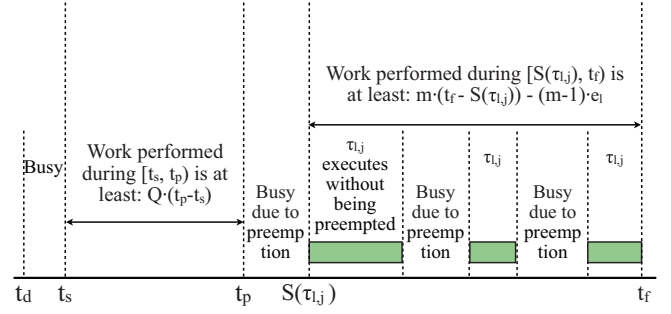


Fig. 7. Subcase 2.3.

Since W is defined to be the amount of work due to jobs in \mathbf{d} that can compete with τ_{lj} after $t_d + y$ and $W \leq 0$, the latest completion time of τ_{lj} is at $t_d + y + e_l - \eta_l$. Therefore, the response time of τ_{lj} is

$$\begin{aligned}
t_f - r_{lj} &\leq t_d + y + e_l - \eta_l - r_{lj} \\
&= y + e_l - \eta_l + (t_d - r_{lj}) \\
&= y + e_l - \eta_l + p_l \\
&\{ \text{by(11)} \} \\
&= \frac{Q}{m} \cdot x + \frac{\eta_l}{m} + e_l - \eta_l + p_l \\
&\{ \text{by(9)} \} \\
&\leq x + e_l + p_l.
\end{aligned}$$

Subcase 2.3. $t_p > t_s$. The earliest time τ_{lj} can commence execution is t_p , as shown in Fig. 7. Let $S(\tau_{lj})$ be the time when τ_{lj} starts execution for the first time. If τ_{lj} is not preempted after t_p , then τ_{lj} starts execution at t_p and completes no later than $t_p + e_l^{\min}$. Thus, we have

$$\begin{aligned}
t_f - r_{lj} &= t_p + e_l^{\min} - r_{lj} \\
&\{ \text{by(4)} \} \\
&\leq t_d - p_l + x + e_l + e_l^{\min} - r_{lj} \\
&= x + e_l + e_l^{\min} \\
&\{ \text{because } e_l^{\min} \leq p_l \} \\
&\leq x + e_l + p_l.
\end{aligned}$$

The other possibility is that τ_{lj} gets preempted after t_p . Let λ denote the set of tasks including τ_l that have ready jobs in \mathbf{d} at any time instant within $[t_s, t_p)$.

We now prove that $|\lambda| \geq Q$ holds. By Definition 8, in order for τ_{lj} to be preempted after t_p , the number of processors required by tasks in λ (note that $\tau_l \in \lambda$) at some time instant after t_p must exceed m . Thus, the maximum total number of threads of tasks in λ that can execute in parallel at the same time must exceed m , which gives

$$\sum_{\tau_l \in \lambda} v_i^{\max} > m. \quad (12)$$

Thus, by the definition of v_{\max_k} , we have $\sum_{k=1}^{|\lambda|} v_{\max_k} \geq \sum_{\tau_l \in \lambda} v_i^{\max} > m$. By Definition 10, we consider two cases: $v_{\max_1} \leq m$ and $v_{\max_1} > m$. If $v_{\max_1} \leq m$, then $|\lambda| \geq Q$ holds. On the other hand, if $v_{\max_1} > m$, then although $\sum_{k=1}^{|\lambda|} v_{\max_k} > m$ may hold when $|\lambda| = 1$, λ clearly needs to contain at least two tasks in order for τ_{lj} to be preempted (namely, τ_l and at least one other task). Thus, $|\lambda| \geq Q$ also holds in this case.

Because $|\lambda| \geq Q$, we know that at least Q tasks have ready jobs in \mathbf{d} at any time instant within $[t_s, t_p)$, which occupy at least Q processors throughout the interval $[t_s, t_p)$. Thus, the amount of work due to \mathbf{d} performed in $[t_s, t_p)$ is at least $Q \cdot (t_p - t_s)$. We now complete the proof of Subcase 3.2 (and thereby Lemma 3).

By the definitions of t_s and t_p , $[t_d, t_s)$ and $[t_p, S(\tau_{l,j}))$ are busy for \mathbf{d} . As discussed above, the amount of work due to \mathbf{d} performed in $[t_s, t_p)$ is at least $Q \cdot (t_p - t_s)$. Moreover, the amount of work due to \mathbf{d} performed in $[S(\tau_{l,j}), t_f)$ is at least $m \cdot (t_f - S(\tau_{l,j})) - (m-1) \cdot e_l$.⁴ Thus, we have

$$LAG(\mathbf{d}, t_d, S) \geq m \cdot (t_s - t_d) + Q \cdot (t_p - t_s) + m \cdot (S(\tau_{l,j}) - t_p) + m \cdot (t_f - S(\tau_{l,j})) - (m-1) \cdot e_l.$$

By (10), we therefore have

$$Q \cdot x - (m-1) \cdot e_l \geq m \cdot (t_s - t_d) + Q \cdot (t_p - t_s) + m \cdot (S(\tau_{l,j}) - t_p) + m \cdot (t_f - S(\tau_{l,j})) - (m-1) \cdot e_l,$$

which gives,

$$t_f - t_d \leq \frac{Q}{m} \cdot x + \left(1 - \frac{Q}{m}\right) \cdot (t_p - t_s). \quad (13)$$

Also, we have $t_p - t_s \leq t_p - t_d \leq t_d - p_l + x + e_l - t_d = x - p_l + e_l$. Therefore,

$$\begin{aligned} t_f - r_{l,j} &= t_f - t_d + p_l \\ &\stackrel{\{\text{by (13)}\}}{\leq} \frac{Q}{m} \cdot x + \left(1 - \frac{Q}{m}\right) \cdot (x - p_l + e_l) + p_l \\ &\stackrel{\{\text{by (9)}\}}{\leq} x + p_l + e_l. \quad \square \end{aligned}$$

3.3. Determining x

Setting the upper bound on $LAG(\mathbf{d}, t_d, S)$ in Lemma 2 to be at most the lower bound in Lemma 3 will ensure that the response time of $\tau_{l,j}$ is at most $x + p_l + e_l$. By solving for the minimum x that satisfies the resulting inequality, we obtain a value of x that is sufficient for ensuring a response time of at most $x + p_l + e_l$. By Lemmas 2 and 3, this inequality is

$$U \cdot x + E \leq Q \cdot x - (m-1) \cdot e_l.$$

Solving for x , we have

$$x \geq \frac{E + (m-1) \cdot e_l}{Q - U}. \quad (14)$$

If x equals the right-hand side of (14), then the response time of $\tau_{l,j}$ will not exceed $x + p_l + e_l$. A value for x that is independent of the parameters of τ_l can be obtained by replacing $(m-1) \cdot e_l$ with $\max_i((m-1) \cdot e_l)$ in (14).

Theorem 1. With x as defined above, the response time for any task τ_l scheduled under GEPPF is at most $x + p_l + e_l$, provided $U < Q$, where U and Q are defined in Definitions 7 and 10, respectively.

3.4. A case with no utilization loss

The following corollary shows that GEPPF results in no utilization loss for scheduling any parallel task system on two processors.

Corollary 1. For two-processor systems, the response time for any task τ_l scheduled under GEPPF is at most $x + p_l + e_l$, where

⁴ We apply the same reasoning as used in Subcase 2.2. All intervals in $[S(\tau_{l,j}), t_f)$ during which $\tau_{l,j}$ is preempted are busy, and $\tau_{l,j}$ can execute for at most e_l time in $[S(\tau_{l,j}), t_f)$. (Within such intervals, at least one processor is occupied by $\tau_{l,j}$.)

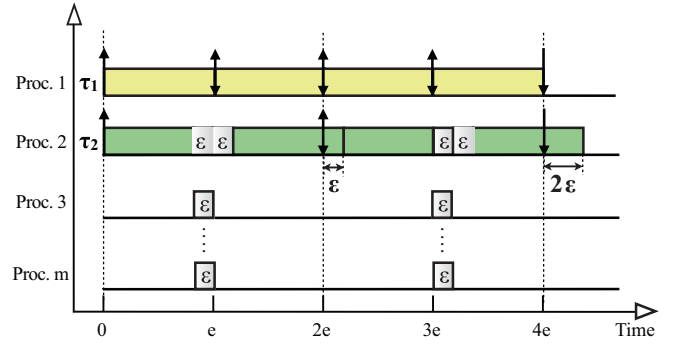


Fig. 8. The worst-case parallel task set.

$x = \frac{E + (m-1) \cdot e_l}{Q - \max_i(u_i)}$ and $\max_i(u_i)$ is the maximum task utilization of tasks in τ .

Proof. If the system only contains one task, then clearly this task, denoted τ_1 , has bounded response time, which is given by $e_1^{min} \leq x + p_1 + e_l$. If the system contains more than one task, then by Definitions 7 and 10 and $m = 2$, we have $U = \max_i(u_i)$ and $Q = 2 = m$. Thus, the utilization constraint in Theorem 1 becomes $\max_i(u_i) < Q = m$, which always holds. \square

3.5. Cases with utilization loss

As shown in Theorem 1 and Corollary 1, the utilization constraint $U < Q$ is needed on $m \geq 3$ processors while no utilization constraint is needed on $m = 2$ processors. By Definitions 7 and 10, in the worst case, $U = U_{sum}$ and $Q = 2$. This implies that in some cases even when m is arbitrarily large, $U_{sum} < 2$ is needed in our analysis. Since no utilization loss can be achieved on two processors as shown in Corollary 1, we can schedule any parallel task system with $U_{sum} = 2$ on only two processors (i.e., leave the other $m-2$ processors idle if $m > 2$). Thus, in the worst case, $U_{sum} \leq 2$ (rather than $U_{sum} < 2$) is needed under our analysis for any parallel task system to have bounded response times for $m \geq 3$ processors. To discern how severe such constraints must fundamentally be, we next show that for any $m \geq 3$, there exists a parallel task system with a total utilization of $2 + \sigma$ that has unbounded response times, where σ can be an arbitrarily small value. This proves that utilization constraints are fundamental for parallel task systems scheduled on $m \geq 3$ processors. (Note that this task set also violates our derived utilization constraint.)

Worst-case parallel task set. Consider a parallel task system containing two parallel tasks. Task τ_1 has only one segment that contains one thread with an execution cost of e time units, and τ_1 has a period of e time units. Thus, τ_1 has a utilization of 1.0. Task τ_2 has three segments, where the first segment contains one thread with an execution cost of $e - \epsilon$ time units, where ϵ can be an arbitrarily small value, the second segment contains m parallel threads, each of which has an execution cost of ϵ time units, and the third segment contains one thread with an execution cost of e time units. τ_2 has a period of $2e$ and a utilization of $\frac{e - \epsilon + m \cdot \epsilon + e}{2e} = 1 + \frac{(m-1)}{2e} \cdot \epsilon$. Thus, this task set has a total utilization of $2 + \frac{(m-1)}{2e} \cdot \epsilon$, or rather $2 + \sigma$, where $\sigma = \frac{(m-1)}{2e} \cdot \epsilon$ can be arbitrarily small.

Fig. 8 shows the GEPPF schedule of this parallel task system on any $m \geq 3$ processors. It is clearly seen that task τ_2 's response time grows unboundedly regardless of m .

4. Optimization

The utilization loss seen in the utilization constraint $U < Q$ is mainly caused by a small value of Q . (Note that by Definition 7, U is completely determined by the tasks' execution costs and periods, which are fixed parameters.) By Definition 10, Q depends on v_{\max_i} ($1 \leq i \leq n$). If the value of v_{\max_i} can be decreased, then the value of Q is increased.

To decrease v_{\max_i} ($1 \leq i \leq n$), we can seek to decrease v_k^{\max} (the maximum number of threads in any segment of τ_k) for each task $\tau_k \in \tau$. This can be done by splitting any segment of τ_k with a large number of threads into multiple sequential sub-segments, each of which has a smaller number of threads, thus decreasing v_k^{\max} . Notice that a critical constraint to enable such splittings is to ensure that $e_k^{\min} \leq p_k$ still holds for any task τ_k after splitting; otherwise, response times may grow unboundedly. Thus, for each task, we

need to determine the maximum degree to which its segments can be split.

We propose algorithm *Q-Optimization* to increase Q for any given parallel task system τ by decreasing v_k^{\max} for each task $\tau_k \in \tau$, as discussed above. The pseudo-code for this algorithm is given in Figs. 9–11. Applying this algorithm can also reduce response time bounds, as seen in Section 5.

Algorithm description. Algorithm *Q-Optimization* seeks to increase the value of Q by decreasing the maximum number of threads in any segment of each task. In the code, v_i^{\max} (v_i^{secmax}) denotes the number of threads in the segment of τ_i with the largest (second largest) number of threads. Note that if all segments of task τ_i contain the same number of threads, then $v_i^{\text{secmax}} = 0$.

We first describe the function *SPLIT* (shown in Fig. 10) used in the main algorithm (shown in Fig. 9). *SPLIT*(τ_k, H) splits the segments with the maximum number of threads into a number of

ALGORITHM: Q-OPTIMIZATION

further-split-flag: BOOLEAN

C_k : INTEGER, INITIALLY $v_k^{\text{secmax}} + 1$

v_k^{\max} : INTEGER, DEFINED IN SECTION 4

v_k^{secmax} : INTEGER, DEFINED IN SECTION 4

h : INTEGER, INITIALLY $h := 1$

A_k : THE SET OF SEGMENTS IN τ_k THAT HAVE v_k^{\max} THREADS

```

1  for each parallel task  $\tau_k \in \tau$ 
2    further-split-flag := false
3    do
4      SPLIT( $\tau_k, v_k^{\text{secmax}}$ )
5      if  $e_k^{\min} < p_k$ 
6        if  $v_k^{\max} \neq 1$ 
7          then further-split-flag := true
8      else if  $e_k^{\min} > p_k$ 
9        Restore the structure of  $\tau_k$  to the one before the last splitting and
10       update segment notations,  $A_k$ ,  $v_k^{\max}$ ,  $v_k^{\text{secmax}}$ , and  $C_k$  accordingly
11       while  $C_k < v_k^{\max}$ 
12         SPLIT( $\tau_k, C_k$ )
13         if  $e_k^{\min} \leq p_k$ 
14           break
15       else
16         Restore the structure of  $\tau_k$  to the one before the last splitting and
17         update segment notations,  $A_k$ ,  $v_k^{\max}$ ,  $v_k^{\text{secmax}}$ , and  $C_k$  accordingly
18          $C_k := C_k + 1$ 
19   while further-split-flag = true

```

Fig. 9. Algorithm *Q-Optimization*.

FUNCTION: SPLIT(τ_k, H)

- 1 **for** each segment $\tau_k^j \in A_k$
- 2 Split τ_k^j into $\lceil \frac{v_k^j}{H} \rceil$ sequential sub-segments, each with at most H threads, and
- 3 assign threads to each sub-segment in smallest-thread-ID-first order and
- 4 update segment notations, A_k , v_k^{max} , v_k^{secmax} , and C_k accordingly
- 5 **COMBINE**(τ_k)
- 6 Calculate e_k^{min}

Fig. 10. Function SPLIT.

FUNCTION: COMBINE(τ_k)

- 1 **while** τ_k^h exists
- 2 **if** τ_k^h and τ_k^{h+1} (if any) are sub-segments that originally belong to the same segment, and $v_k^h + v_k^{h+1} \leq v_k^{max}$
- 3 **then** combine τ_k^h and τ_k^{h+1} into one segment and
- 4 update segment notations, A_k , v_k^{max} , v_k^{secmax} , and C_k accordingly
- 5 **else**
- 6 $h := h + 1$

Fig. 11. Function Combine.

sequential sub-segments, each with at most H threads (lines 1–4 in function SPLIT). (Note that several variables used in this function are defined in algorithm Q-Optimization shown in Fig. 9.) Threads are assigned to each of these sub-segments in smallest-thread-ID-first order, until either a sub-segment contains H threads or all threads have been assigned. Then in line 5, function COMBINE (shown in Fig. 11) seeks to combine any two sub-segments that originally belong to the same segment into one segment if the sum of the number of threads in both sub-segments is no greater than the maximum number of threads of any segment. Finally, function SPLIT calculates e_k^{min} (line 6 in function SPLIT) using the method we discussed in Section 2.

Now we describe algorithm Q-Optimization in detail. First we make two important observations. (i) For any task τ_i that contains at least two segments having a different number of threads, we desire to reduce the number of threads of its segments that contain the maximum number of threads among all segments of τ_i to no less than v_i^{secmax} . Further reductions do not reduce v_i^{max} . (ii) For any task τ_i containing at least two segments that have the same number of threads, we desire to reduce the number of threads of such segments by the same amount. Reducing any such segment's thread count by a greater amount than the others does not reduce v_i^{max} .

Motivated by these two observations, the algorithm first executes SPLIT(τ_k, v_k^{secmax}), which splits each of the segments in τ_k that have the maximum number of threads into a sequential number of sub-segments, each with at most v_k^{secmax} threads. After such a splitting, if $e_k^{min} < p_k$ and $v_k^{max} \neq 1$ (lines 5–7 in algorithm Q-Optimization), then we set the *further-split-flag* to be true, which implies that there is still the potential for us to split τ_k to further reduce v_k^{max} .

On the other hand, if $e_k^{min} > p_k$ after such a splitting (line 8 in algorithm Q-Optimization), then it implies that such a splitting causes e_i^{min} to exceed τ_k 's period (which causes τ_k to have un-

bounded response times) and is thus invalid. Since this splitting is invalid, we restore the task structure to the one before the splitting (lines 9–10 in algorithm Q-Optimization). Thus, we now know that it is impossible to split segments in τ_k to reduce v_k^{max} to equal v_k^{secmax} . However, by splitting, we might still be able to reduce v_k^{max} to some number between v_k^{secmax} and v_k^{max} (realized by lines 12–14 in algorithm Q-Optimization). Note that the minimum value of such a number is given by C_k (for otherwise it would have been possible to reduce v_k^{max} to equal v_k^{secmax} given that C_k is initially $v_k^{secmax} + 1$). Therefore, starting from C_k , the algorithm uses the SPLIT function and compares the resulting e_k^{min} with p_k to determine whether any such splitting is valid (lines 11–18 in algorithm Q-Optimization using the logic discussed above).

Optimization example. Since we seek to decrease v_k^{max} for each task τ_k in any given task system using the same optimization algorithm, we use one example task τ_1 to illustrate the idea. In this example, $m = 4$ and τ_1 originally has five segments, as illustrated in Fig. 12(a). The notation $\tau_1^{ij}(e)$ in Fig. 12 denotes that thread τ_1^{ij} has an execution cost of e time units. τ_1 has a period of 18 time units, thus $p_1 = 18$.

Because we want to decrease v_1^{max} , we first try to decrease the number of threads of segments in τ_1 that have the largest number of threads, which are τ_1^2 and τ_1^3 (realized by executing algorithm Q-Optimization). Therefore, according to observations (i) and (ii) discussed above, we split each of τ_1^2 and τ_1^3 into two sequential sub-segments, one with three threads and the other one with one thread (realized by executing line 4 in algorithm Q-Optimization), as shown in Fig. 12(b) (note that in the figure updated segment notations are used after each splitting). After this splitting, we obtain $e_1^{min} = 15 < p_1 = 18$ (we apply the same method discussed in Section 2 to obtain e_1^{min}). Thus, this splitting is valid (as verified in lines 5–7 in algorithm Q-Optimization). Now we obtain a task τ_1 in which segments τ_1^2 , τ_1^4 , and τ_1^6 have the largest number of threads (three threads per segment), while segment τ_1^1 has the

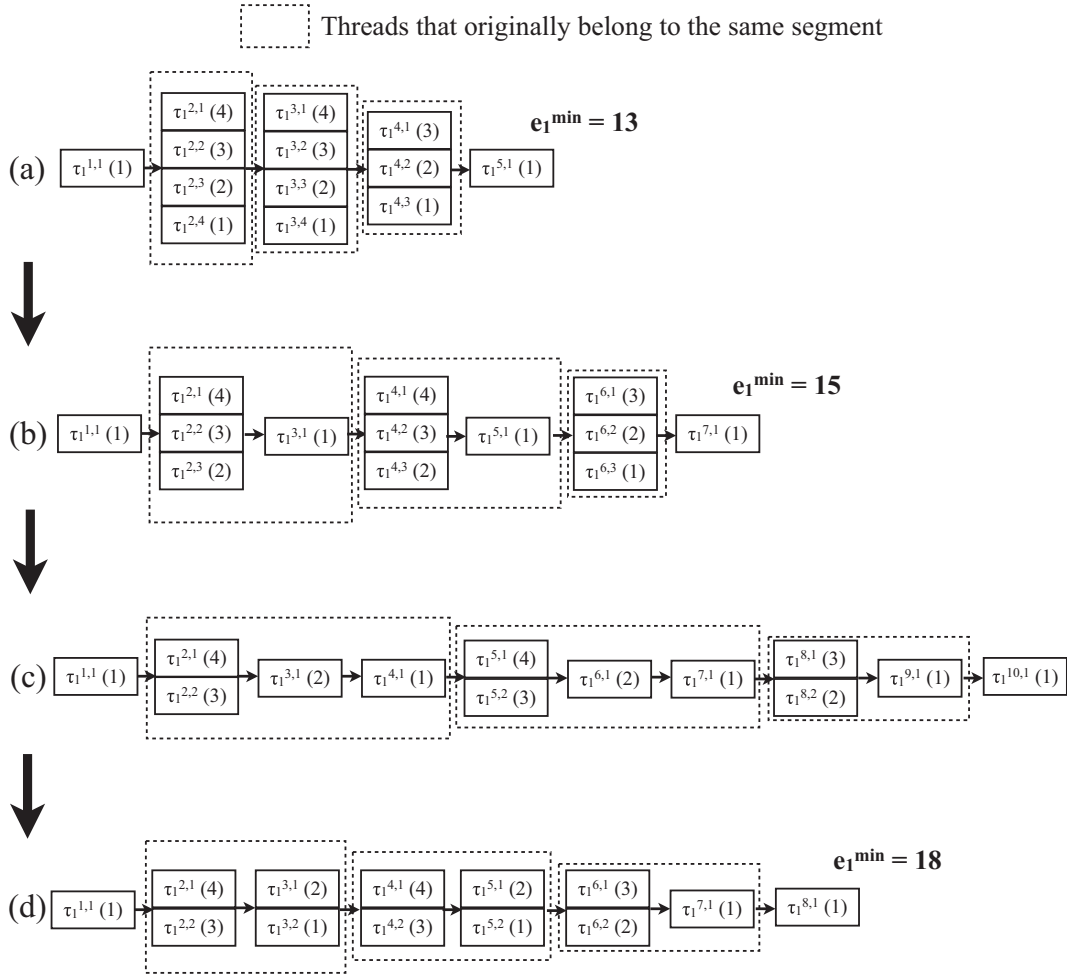


Fig. 12. Illustration of the optimization algorithm.

second largest number of threads (one thread per segment). Therefore, we again try to reduce the number of threads of τ_1^2 , τ_1^4 , and τ_1^6 to no less than the number of threads of τ_1^1 . This can be achieved by splitting each of these three segments into three sequential segments, each of which contains only one thread (again, realized by executing line 4 in algorithm *Q-Optimization*). However, after such a splitting, we have $e_1^{\min} = 28 > p_1 = 18$. Thus, such a splitting is invalid (as verified in lines 8–10 in algorithm *Q-Optimization*).

Therefore, our goal now is trying to reduce the number of threads of τ_1^2 , τ_1^4 , and τ_1^6 to a smallest possible number, which is two threads per segment in this case (realized by executing lines 11–18 in algorithm *Q-Optimization*). As shown in Fig. 12(c), we split each of τ_1^2 , τ_1^4 , and τ_1^6 into two sequential sub-segments, one with two threads and another one with one thread. Also notice that after this splitting, τ_1^3 and τ_1^4 originally belonged to the same segment, and τ_1^5 and τ_1^7 originally belonged to the same segment. Since combining τ_1^3 and τ_1^4 (as well as τ_1^5 and τ_1^7) into one sub-segment does not increase v_1^{\max} , we combine them in such a way to decrease e_1^{\min} (realized by executing function *COMBINE*), as illustrated in Fig. 12(d). After this splitting, we have $e_1^{\min} = 18 = p_1$. Thus, we cannot split segments any further (as verified in lines 13–14 in algorithm *Q-Optimization*), and we successfully reduce v_1^{\max} from 4 to 2.

5. Experimental evaluation

In this section, we describe experiments conducted using randomly-generated parallel task sets to evaluate the applicability of

the response time bound in Theorem 1. Moreover, we evaluate whether the optimization algorithm can effectively improve schedulability (with respect to bounded response times) and reduce the bound.

Experimental setup. In our experiments, parallel task sets were generated as follows. The number of segments of each task was uniformly distributed over [1,30]. The number of threads of each segment was distributed differently for each experiment using three uniform distributions: $[1, m/2]$ (*low parallelism*), $[m/2, m]$ (*high parallelism*), and $[1, m]$ (*random parallelism*), where m is the number of processors. The execution cost of each thread was uniformly distributed over $[1\text{ms}, 100\text{ms}]$. The worst-case execution cost e_i and the best-case execution cost e_i^{\min} of each parallel task τ_i were then calculated using the approach discussed in Section 2. Then, for each task τ_i , its period was uniformly distributed over $[e_i^{\min}, e_i^{\min} + e_i]$, and its utilization was calculated using e_i and p_i . We also varied the system utilization U_{sum} within $\{0.1, 0.2, \dots, m\}$. For each U_{sum} , 1000 parallel task sets were generated for systems with four, six, and eight processors.⁵ Each such parallel task set was generated by creating parallel tasks until total utilization exceeded U_{sum} , and by then reducing the last task's utilization so that the total system utilization equaled U_{sum} . For each generated system, we first checked schedulability (i.e., the ability to

⁵ For systems with higher processor counts, recent experimental work [2] suggests that when overheads are considered, clustered scheduling approaches (where groups of processors with low processor counts that share low-level caches are scheduled globally) are better than global approaches.

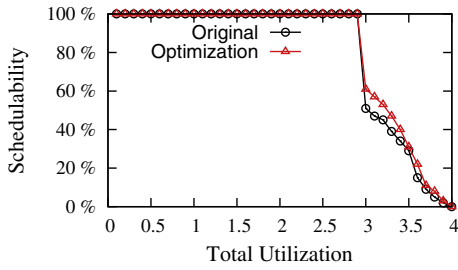


Fig. 13. Schedulability: $m = 4$, low parallelism.

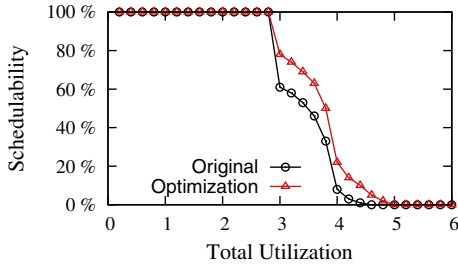


Fig. 14. Schedulability: $m = 6$, low parallelism.

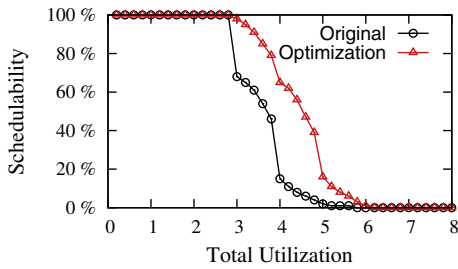


Fig. 15. Schedulability: $m = 8$, low parallelism.

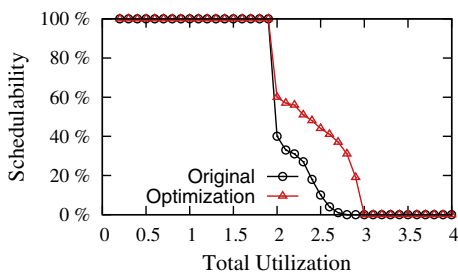


Fig. 16. Schedulability: $m = 4$, high parallelism.

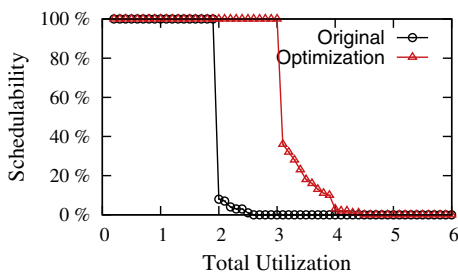


Fig. 17. Schedulability: $m = 6$, high parallelism.

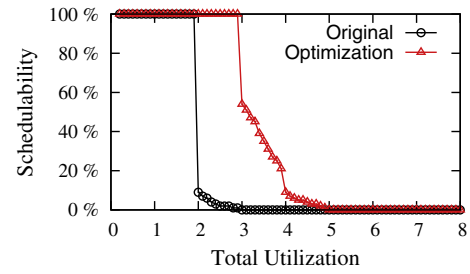


Fig. 18. Schedulability: $m = 8$, high parallelism.

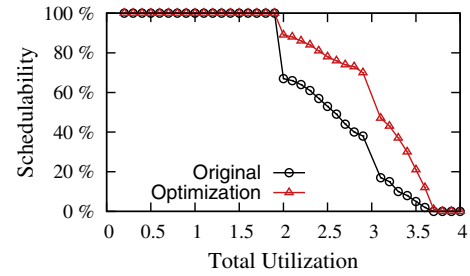


Fig. 19. Schedulability: $m = 4$, random parallelism.

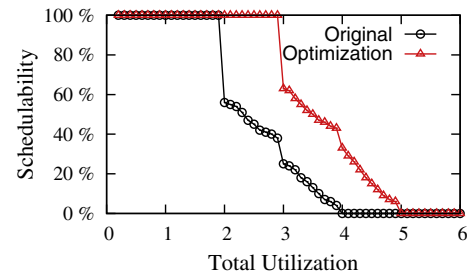


Fig. 20. Schedulability: $m = 6$, random parallelism.

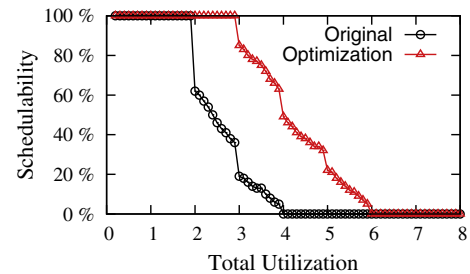


Fig. 21. Schedulability: $m = 8$, random parallelism.

ensure bounded response times) and the magnitude of response time bounds using [Theorem 1](#). Then, for each such generated system, we applied the optimization algorithm and re-checked schedulability and response time bounds. In doing so, system overheads were ignored (factoring overheads into our analysis is beyond the scope of this paper). In all figures and tables presented in this section, we let “Original” and “Optimization” denote results under the original analysis and results after applying the optimization algorithm Q-Optimization.

Results. The schedulability results that were obtained on four-, six-, and eight-processor systems with different degrees of intra-task parallelism are shown in [Figs. 13–21](#), respectively. In these

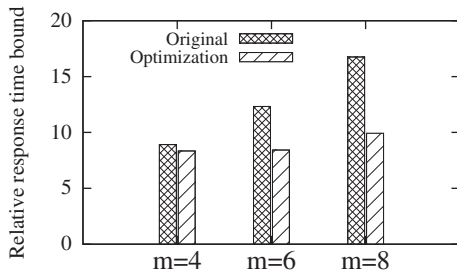


Fig. 22. Response time bounds: low parallelism.

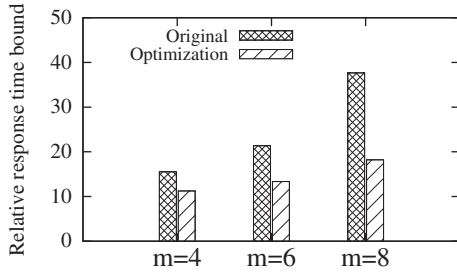


Fig. 23. Response time bounds: high parallelism.

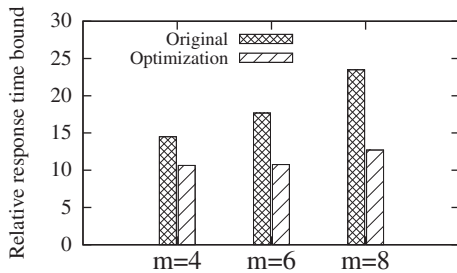


Fig. 24. Response time bounds: random parallelism.

figures, the x -axis denotes U_{sum} and each curve plots the fraction of the generated parallel task sets the corresponding approach successfully scheduled, as a function of U_{sum} . As seen, our analysis can provide reasonable schedulability. For example, as shown in Fig. 13, on four processors with low parallelism, all parallel task sets have bounded response times until U_{sum} reaches 3.0 and more than 40% of the task sets still have bounded response times when U_{sum} reaches 3.3. Moreover, the optimization algorithm is able to effectively improve schedulability, especially when the processor count is large or the intra-task parallelism is high. For example, as illustrated in Fig. 21, on eight processors with random parallelism, the optimization algorithm can improve schedulability by more than 400% in many cases (e.g., when $U_{sum} = 3.0$). Such improvements tend to increase with increasing processor count or increasing parallelism. This is because when m becomes larger or the number of threads per segment increases, it is easier to increase Q by applying the optimization algorithm, which is intuitive according to the definition of Q . Note that, when schedulability drops significantly, it does so at an integral values of U_{sum} . For example, as seen in Fig. 13, when U_{sum} reaches 3.0, schedulability drops from 100% to less than 50% under Original. This is because when U_{sum} reaches 3.0, by Definition 7, U may also equal 3.0 since some parallel tasks very likely have utilization greater than 1.0. Thus, Q has to be 4.0 instead of 3.0 (when the utilization is below 3.0) in order for the utilization constraint $Q > U$ to hold; this obviously makes this constraint much more severe.

Figs. 22–24 show the computed response time bounds using Theorem 1 under Original and Optimization. To better illustrate the magnitude of the response time bounds, we plot *relative response time bounds*. A task's relative response time bound is given by the ratio of its response time bound divided by its period. The data in Fig. 13 shows average relative response time bounds obtained by considering all tasks in certain selected task sets. Such task sets were selected by considering values of U_{sum} for which 100% schedulability can be ensured, which guarantees all such task sets valid response time bounds. For example, on four processors, we calculated the average relative response time bound over task sets whose utilizations are within $[0.1, 3)$ (all such task sets are schedulable and thus have valid response time bounds). As seen in the figure, our analysis can achieve reasonable response time bounds. For example, as shown in Fig. 22, on four processors with low parallelism, the average relative response time bound is around nine. The benefit of the optimization algorithm is apparent. For example, as illustrated in Fig. 23, on eight processors with high parallelism, we can reduce the average relative response time bound from around 33 to less than 18. This is because applying the optimization algorithm only increases Q and does not change other values in the response time bound expression shown in Theorem 1.

6. Conclusion

We have presented schedulability analysis for sporadic parallel task systems under GEPPF scheduling. The proposed analysis shows that such systems can be efficiently supported on multiprocessors with bounded response times. In experiments presented herein, our analysis is proved to provide good performance with respect to both schedulability and response time bounds. In future work, it would be interesting to investigate more practical parallel task models where data is communicated among segments within a parallel task. Moreover, allowing more general parallel execution patterns such as cycles would expand the applicability of our results.

References

- [1] D. Bailey, An optimal scheduling algorithm for parallel video processing, in: Proc. of the 5th IEEE Intl. Conf. on Multimedia Computing and Systems, 1998, pp. 245–248.
- [2] A. Bastoni, B. Brandenburg, J. Anderson, An empirical comparison of global, partitioned, and clustered multiprocessor real-time schedulers, in: Proc. of the 31st Real-Time Sys. Symp., 2010, pp. 14–24.
- [3] R. Chandra, R. Menon, L. Dagum, D. Kohn, D. Maydan, J. McDonald, *Parallel Programming in OpenMP*, Morgan Kaufman, 2000.
- [4] S. Chen, S. Schlosser, MapReduce meets wider varieties of applications, Technical Report IRP-TR-08-05, Intel Labs Pittsburgh, 2008.
- [5] J. Dean, S. Ghemawat, MapReduce: simplified data processing on large clusters, in: Proc. of the 6th USENIX Conf. on Symp. on Operating System Design and Implementation, 2004, pp. 137–150.
- [6] U. Devi, J. Anderson, Tardiness bounds under global EDF scheduling on a multiprocessor, in: Proc. of the 26th IEEE Intl Real-Time Sys. Symp., 2005, pp. 330–341.
- [7] E. Dougherty, P. Laplante, *Introduction to Real-time Imaging*, Wiley-IEEE Press, 1995.
- [8] D. Feitelson, Job scheduling in multiprogrammed parallel systems, Technical Report RC 19790 (87657), 1997.
- [9] D. Feitelson, L. Rudolph, Parallel job scheduling: issues and approaches, in: Proc. of the Workshop on Job Scheduling Strategies for Parallel Processing, 1995, pp. 1–18.
- [10] E. Horowitz, S. Sahni, Exact and approximate algorithms for scheduling nonidentical processors, *J. ACM* 23 (2) (1976).
- [11] Y. Kitamura, A. Smith, H. Takemura, F. Kishino, Parallel algorithms for real-time colliding face detection, in: Proc. of the 4th IEEE Workshop on Robot and Human Communication, 1995, pp. 211–218.
- [12] K. Lakshmanan, S. Kato, R. Rajkumar, Scheduling parallel real-time tasks on multi-core processors, in: Proc. of the 31st Real-Time Sys. Symp., 2010, pp. 259–268.

- [13] H. Leontyev, J. Anderson, Generalized tardiness bounds for global multiprocessor scheduling, in: Proc. of the 28th IEEE Real-Time Systems, 2007, pp. 413–422.
- [14] C. Liu, The real-time multi-resource task model, in: Proc. of the 3rd Intl. Real-Time Scheduling Open Problems Seminar, 2012.
- [15] C. Liu, J. Anderson, Supporting soft real-time parallel applications on multicore processors, in: Proc. of the 18th IEEE International Conference on Embedded and Real-Time Computing Systems and Applications, 2012, pp. 114–123.
- [16] C. Liu, J. Anderson, Improving the schedulability of sporadic self-suspending soft real-time multiprocessor task systems, in: Proc. of the 16th IEEE Int'l Conf. on Embedded and Real-Time Computing Sys. and Apps., 2010, pp. 14–23.
- [17] C. Liu, J. Anderson, Supporting graph-based real-time applications in distributed systems, in: Proc. of the 17th IEEE Int'l Conf. on Embedded and Real-Time Computing Sys. and Apps., 2011, pp. 143–152.
- [18] C. Liu, J. Anderson, Scheduling suspendable, pipelined tasks with non-preemptive sections in soft real-time multiprocessor systems, in: Proc. of the 16th IEEE Real-Time and Embedded Tech. and Apps. Symp., 2010, pp. 23–32.
- [19] C. Liu, J. Anderson, Supporting pipelines in soft real-time multiprocessor systems, in: Proc. of the 21st Euromicro Conf. on Real-Time Sys., 2009, pp. 269–278.
- [20] C. Liu, J. Anderson, Supporting sporadic pipelined tasks with early-releasing in soft real-time multiprocessor systems, in: Proc. of the 15th IEEE Int'l Conf. on Embedded and Real-Time Computing Sys. and Apps., 2009, pp. 284–293.
- [21] C. Liu, J. Anderson, A new technique for analyzing soft real-time self-suspending task systems, *ACM SIGBED Rev.* (2012) 29–32.
- [22] C. Liu, J. Anderson, Supporting soft real-time dag-based systems on multiprocessors with no utilization loss, in: Proc. of the 31st Real-Time Sys. Symp., 2010, pp. 3–13.
- [23] C. Liu, J. Anderson, An $O(m)$ analysis technique for supporting real-time self-suspending task systems, in: Proc. of the 33th IEEE Real-Time Systems Symposium, 2012, pp. 373–382.
- [24] C. Liu, J. Anderson, Task scheduling with self-suspensions in soft real-time multiprocessor systems, in: Proc. of the 30th Real-Time Sys. Symp., 2009, pp. 425–436.
- [25] C. Liu, J. Anderson, Suspension-Aware Analysis for Hard Real-Time Multiprocessor Scheduling, in: Proc. of the 25th Euromicro Conference on Real-Time Systems, 2013, pp. 271–281.
- [26] C. Liu, J. Li, W. Huang, J. Rubio, E. Speight, X. Lin, Power-efficient time-sensitive mapping in cpu/gpu heterogeneous systems, in: Proc. of the 21st International Conference on Parallel Architectures and Compilation Techniques, 2012, pp. 23–32.
- [27] G. Nelissen, V. Berten, J. Goossens, D. Milojevic, Techniques optimizing the number of processors to schedule multi-threaded tasks, in: Proc. of the 24th Euromicro Conf. on Real-Time Sys., 2012, pp. 321–330.
- [28] G. Nelissen, V. Berten, V. Nelis, J. Goossens, D. Milojevic, An unfair but optimal multiprocessor scheduling algorithm for sporadic tasks, in: Proc. of the 24th Euromicro Conf. in Real-Time Sys., 2012.
- [29] A. Saifullah, K. Agrawal, C. Lu, C. Gill, Multi-core real-time scheduling for generalized parallel task models, in: Proc. of the 32nd Real-Time Sys. Symp., 2011, pp. 217–226.
- [30] A. Srinivasan, J. Anderson, Fair scheduling of dynamic task systems on multiprocessors, *J. Syst. Softw.* 77 (1) (2005) 67–80.
- [31] C. Volker, V. Hamscher, R. Yahyapour, Economic scheduling in grid computing, in: Proc. of the Conf. on Scheduling Strategies for Parallel Processing, 2002, pp. 128–152.
- [32] M. Zaharia, D. Borthakur, J. Sarma, K. Elmeleegy, S. Shenker, I. Stoica, Job scheduling for multi-user MapReduce clusters, Technical Report EECS-2009055, UC Berkeley, 2009.



Cong Liu is an assistant professor in the Department of Computer Science at the University of Texas at Dallas. He received his Ph.D. degree in Computer Science from the University of North Carolina at Chapel Hill in July 2013. His research interests include real-time and embedded systems, operating systems, cloud computing, and wireless sensor networks. He has published over 30 papers in premier conferences and journals such as RTSS, PACT, ECRTS, RTAS, and JPDC. He received the Best Student Paper Award at the 30th IEEE Real-Time Systems Symposium, the premier real-time and embedded systems conference. He also received the best papers award at the 17th RTCSA and the UNC's dissertation fellowship.



James H. Anderson is a professor in the Department of Computer Science at the University of North Carolina at Chapel Hill. He received a B.S. degree in Computer Science from Michigan State University in 1982, an M.S. degree in Computer Science from Purdue University in 1983, and a Ph.D. degree in Computer Sciences from the University of Texas at Austin in 1990. In 1995, Dr. Anderson received the U.S. Army Research Office Young Investigator Award, and in 1996, he was named Alfred P. Sloan Research Fellow. He has served as program chair and/or general chair for several conferences, including the ACM Symposium on Principles of Distributed Computing, the IEEE International Real-Time Systems Symposium, and the International Conference on Principles of Distributed Systems.