# An Optimal Semi-Partitioned Scheduler Assuming Arbitrary Affinity Masks\*

Sergey Voronov<sup>1</sup> and James H. Anderson<sup>2</sup>

- 1 Department of Computer Science, University of North Carolina at Chapel Hill rdkl@cs.unc.edu
- 2 Department of Computer Science, University of North Carolina at Chapel Hill anderson@cs.unc.edu

#### — Abstract

Modern operating systems like Linux allow task migrations to be restricted by specifying pertask processor affinity masks. Such a mask specifies the set of processor cores upon which a task can be scheduled. In this paper, a semi-partitioned scheduler, AM-Red (affinity mask reduction), is presented for scheduling implicit-deadline sporadic tasks with arbitrary affinity masks on an identical multiprocessor. AM-Red is obtained by applying an affinity-mask-reduction method that produces affinities in accordance with those specified, without compromising feasibility, but with only a linear number of migrating tasks. It functions by employing a tunable frame F of size |F|. For any choice of |F|, AM-Red is soft-real-time optimal, with tardiness bounded by |F|, but the frequency of task migrations is proportional to |F|. If |F| divides all task periods, then AM-Red is also hard-real-time-optimal (tardiness is zero). AM-Red is the first optimal scheduler proposed for arbitrary affinity masks without future knowledge of all job releases. Experiments are presented that assess its practical viability.

1998 ACM Subject Classification C.1.2 [Processor Architectures]: Multiple Data Stream Architectures (Multiprocessors); C.3 [Special-Purpose and Application-Based Systems]—real-time and embedded systems; D.4.1 [Operating Systems]: Process Management—scheduling; D.4.7 [Operating Systems]: Organization and Design—real-time systems and embedded systems; J.7 [Computer Applications]: Computers in Other Systems—real time

**Keywords and phrases** optimal schedulers, processor affinity masks, semi-partitioned scheduling, tardiness analysis

Digital Object Identifier 10.4230/LIPIcs.ECRTS.2018.YY

## 1 Introduction

On multicore machines, particularly ones with relatively high core counts, it is often desirable to limit task migrations to lessen cache- and I/O-related overheads [10] and to enable load balancing [25], among other reasons [18]. Processor affinity masks are an operating-system (OS) mechanism that enables allowed migrations to be flexibly determined. A given task's affinity mask specifies which cores it is allowed to execute upon. General-purpose OSs that support affinity masks include Windows, Linux, and MacOS X. Real-time OSs (RTOSs) that support them include FreeRTOS [15], QNX [1], VxWorks [2], and many others.

Unfortunately, in the real-time systems domain, no optimal scheduler has heretofore been proposed that allows arbitrary affinity masks. Thus, while OSs provide flexible control over

<sup>\*</sup> Work supported by NSF grants CNS 1409175, CPS 1446631, CNS 1563845, and CNS 1717589, ARO grant W911NF-17-1-0294, and funding from General Motors.

migrations through affinity masks in theory, such support must typically be restricted in practice. For example, under Linux's SCHED\_DEADLINE scheduler [9], affinity masks are essentially ignored: any task is assumed to be executable on any core [3]

In this paper, we show for the first time that the goals of allowing arbitrary affinity masks and scheduling real-time tasks optimally do not fundamentally conflict. We do this by presenting a new scheduler, AM-Red (affinity  $\underline{m}$  ask  $\underline{red}$  uction), that optimally schedules implicit-deadline sporadic task sets. Before delving into notable specifics concerning AM-Red, we first review relevant prior work to provide context.

Related work. The existing literature pertaining to scheduling real-time tasks with affinity masks is not very extensive. For hard real-time (HRT) implicit-deadline sporadic task systems, Baruah et al. [6] proposed an exact feasibility test and corresponding scheduler. However, their scheduler relies on clairvoyant knowledge of all job releases, gives rise to impractically frequent task migrations, and has an offline phase with high time complexity. Muneeswari et al. [17] presented a scheduler supporting affinities that they claimed is applicable to real-time systems, but they provided no analysis to support this claim.

Hierarchical affinity masks are often used in practice to reflect multi-level cache hierarchies (L1, L2, etc.) found in multicore machines. With hierarchical affinities, if the masks of two tasks intersect, one must be contained within the other. For hierarchical affinities, Bonifaci et al. [7] proposed a HRT scheduler that ensures a certain "greedy" property that avoids wasted processing capacity. However, they provided no schedulability test.

Contributions. The main contribution of this paper is AM-Red, a new scheduling algorithm for implicit-deadline sporadic task systems with arbitrary affinity masks. AM-Red is a *semi-partitioned* scheduler; under such schedulers, only certain tasks are allowed to migrate and these tasks are determined in an offline allocation phase [4].

AM-Red schedules tasks by iteratively considering a schedule computed offline for a window of time called a *frame* and denoted F. The frame size |F| is a configurable parameter. For soft real-time (SRT) task systems that require bounded deadline tardiness, AM-Red is optimal and ensures a tardiness bound of |F|. The frame size |F| also determines the frequency of task migrations, so choosing |F| yields a tradeoff: larger values reduce migration costs while smaller values reduce tardiness. If |F| divides all task periods, then AM-Red is also optimal for scheduling HRT task systems. For n tasks executing on m processors, if masks are hierarchical, then AM-Red requires O(m+n) time complexity for its offline phase and O(1) time per scheduling decision; these time bounds are asymptotically optimal. To the best of our knowledge, AM-Red is the first non-clairvoyant scheduler to be proposed that is HRT/SRT-optimal for implicit-deadline sporadic tasks under arbitrary affinity masks.

In addition to presenting AM-Red, we also explore a number of issues concerning affinity-mask reductions. In particular, we consider affinity graphs that aggregate the specified masks of all tasks and present a method that can reduce the number of edges in such a graph without compromising task-set feasibility. While feasibility can be assessed using the test of Baruah and Brandenburg [6], we instead use a test proposed here that has lower time complexity. Our reduction method yields affinity masks under which at most (m-1) tasks migrate. This property is actually instrumental in enabling a semi-partitioned approach.

In order to assess the efficacy of AM-Red, we experimentally compared it to the (non-optimal) scheduler proposed in [7], which assumes hierarchical masks, on the basis of migration frequency and task response times. In these experiments, AM-Red demonstrated performance even better than our analysis predicts.

**Organization.** In the rest of this paper, we provide needed background (Sec. 2), present our new feasibility test (Sec. 3), develop algorithm AM-Red by considering first "loop-free"

affinities (Sec. 4) and then arbitrary ones (Sec. 5), consider the special case of hierarchical affinities (Sec. 6), present our experimental evaluation (Sec. 7), and conclude (Sec. 8).

### 2 Background

We consider the problem of scheduling n implicit-deadline sporadic tasks,  $\tau_1, ..., \tau_n$ , on m identical unit-speed cores,  $\pi_1, ..., \pi_m$ . We assume familiarity with the implicit-deadline sporadic task model, consider only task sets in accordance with this model hereafter, and assume that all time-related parameters are rational. We will use the following notation:  $C_i$  denotes the worst-case execution time of task  $\tau_i$ ,  $T_i$  denotes its period,  $D_i = T_i$  denotes its relative deadline, and  $U_i = C_i/T_i \le 1$  denotes its utilization;  $J_{i,k}$  denotes the  $k^{th}$  job released by  $\tau_i$ , and  $C_{i,k} \le C_i$  denotes the execution time of  $J_{i,k}$ ;  $U = \sum_i U_i$  denotes the total system utilization. Job  $J_{i,k}$  has an absolute deadline  $D_i$  time units after its release, and once it has received a processor allocation equal to  $C_{i,k}$ , it is completed. Job  $J_{i,k+1}$  cannot be scheduled until the prior job of  $\tau_i$ ,  $J_{i,k}$ , has completed, even if  $J_{i,k}$  misses its (absolute) deadline.

If a job has a deadline at time  $t_d$  and completes at time  $t_c$ , then its *tardiness* is defined as  $\max(0, t_c - t_d)$ . The tardiness of task  $\tau_i$  is the supremum of the tardiness of any of its jobs. If this value is finite, then we say that  $\tau_i$  has *bounded* tardiness.

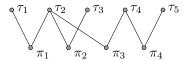
A task set  $\tau$  is HRT-schedulable (resp., SRT-schedulable) under scheduling algorithm S if each task in  $\tau$  has zero (resp., bounded) tardiness in any schedule for  $\tau$  generated by S. A task set  $\tau$  is HRT-feasible (resp., SRT-feasible) if, for any job release sequence (as allowed by the sporadic task model), a schedule exists in which each task has zero (resp., bounded) tardiness. Scheduling algorithm S is HRT-optimal (resp., SRT-optimal) if every HRT-feasible (resp., SRT-feasible) task set  $\tau$  is HRT-schedulable (resp., SRT-schedulable) under S. Although HRT- and SRT-feasibility are fundamentally different concepts in some contexts, we show later that in the context of this paper, they are actually equivalent.

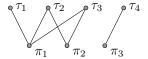
Affinity masks. In practice, affinity masks are usually specified using bit-vectors, but we opt for a more abstract specification. In particular, we define the affinity mask  $\alpha_i$  of task  $\tau_i$  to be the set of cores upon which  $\tau_i$  is allowed to execute. We define the aggregate affinity mask of a subset of tasks  $\tau' \subseteq \tau$  as  $\alpha_{\tau'} = \bigcup_{i \in \tau'} \alpha_i$ . We call the aggregate affinity mask  $\alpha_{\tau}$  of the set of all tasks  $\tau$  the system affinity mask. For a given task set  $\tau$ , we define a bipartite undirected graph called an affinity graph, denoted  $AG(\tau)$ , as follows:  $AG(\tau)$  has n vertices  $\tau_1, \ldots, \tau_n$  (representing tasks), and m vertices  $\pi_1, \ldots, \pi_m$  (representing cores), and contains edge  $(\tau_i, \pi_j)$  if and only if  $\pi_j \in \alpha_i$  (i.e., task  $\tau_i$  can execute on core  $\pi_j$ ).

▶ Example 1. Consider a task set  $\tau$  with five tasks,  $\tau_1, ..., \tau_5$ , to be scheduled on four cores,  $\pi_1, ..., \pi_4$ , with affinity masks  $\alpha_1 = \{1\}$ ,  $\alpha_2 = \{1, 2, 3\}$ ,  $\alpha_3 = \{2\}$ ,  $\alpha_4 = \{3, 4\}$ , and  $\alpha_5 = \{4\}$ .  $AG(\tau)$  is shown in Fig. 1a. The aggregate affinity mask for the subset of tasks  $\{\tau_1, \tau_4, \tau_5\}$  is  $\{\pi_1, \pi_3, \pi_4\}$ , while for  $\{\tau_2, \tau_3\}$  it is  $\{\pi_1, \pi_2, \pi_3\}$ .

Important affinity-mask categories. For a given task set  $\tau$ , we call  $\alpha_{\tau}$  and  $AG(\tau)$  loop-free if and only if  $AG(\tau)$  is acyclic. For example, the affinity graph in Fig. 1a is loop-free. Note that, if  $AG(\tau)$  is loop-free, then, for any i and j, where  $i \neq j$ ,  $|\alpha_i \cap \alpha_j| \leq 1$  holds. (On the other hand, this condition does not necessarily imply that  $AG(\tau)$  is loop-free.)

As mentioned in Sec. 1, hierarchical affinity masks have received prior attention [7, 8, 22]. For a given task set  $\tau$ , we call  $\alpha_{\tau}$  and  $AG(\tau)$  hierarchical if and only if, for any i and j,  $\alpha_i \cap \alpha_j = \emptyset$  or  $\alpha_i \subseteq \alpha_j$  or  $\alpha_i \supseteq \alpha_j$ . Note that hierarchical masks may or may not be loop-free. Note also that the core assignments of any global, clustered, or partitioned scheduler can be





(a) Task set from Ex. 1.

- (b) Task set from Ex. 2 (hierarchical affinities).
- **Figure 1** Example affinity graphs.

$$\tau \text{ is HRT-feasible} \xrightarrow{\text{def}} \tau \text{ is SRT-feasible} \xrightarrow{\text{Lemma 4}} \text{UB}_{\text{6}} \text{Test holds for } \tau$$

$$\text{APA-feas}(\tau, \pi) \text{ [6] holds for } \tau \xrightarrow{\text{Lemma 7 and Theorem 8}} \text{MF Test holds for } \tau$$

**Figure 2** Feasibility proof overview.

specified using masks that are hierarchical. Under global and clustered scheduling, if at least two tasks are allowed to share at least two cores, then such masks will not be loop-free.

If we place no restrictions on affinity masks, then  $\alpha_{\tau}$  and  $AG(\tau)$  are called arbitrary.

▶ Example 2. Consider a task set  $\tau$  with four tasks,  $\tau_1, ..., \tau_4$ , to be scheduled on three cores,  $\pi_1, \pi_2, \pi_3$ , with affinity masks  $\alpha_1 = \{1\}$ ,  $\alpha_2 = \{1, 2\}$ ,  $\alpha_3 = \{1, 2\}$ ,  $\alpha_4 = \{3\}$ . The affinity graph  $AG(\tau)$  for this task set is shown in Fig. 1b. From the figure, it is easy to see that these affinities are hierarchical but not loop-free.

## 3 Feasibility

In order to design a scheduling algorithm for task sets with arbitrary affinity masks and show that it is optimal, we must have a means for determining which such task sets are feasible. A test for HRT-feasibility has been given previously by Baruah *et al.* [6]. In this section, we show that in the considered context, HRT-feasibility and SRT-feasibility are equivalent. We also present a feasibility test that is more efficient than that of Baruah *et al.* 

It is easy to see that HRT-feasibility implies SRT-feasibility because zero tardiness implies bounded tardiness. To show the equivalence of the two, we must therefore show that every SRT-feasible task set  $\tau$  is also HRT-feasible. We do this in three steps: we first establish a special property of SRT-feasible task sets in Sec. 3.1; we then develop a new exact feasibility test in Sec. 3.2 based on max flow; finally, we show the equivalence of the obtained test to that of Baruah et al. [6]. A depiction of these steps is given in Fig. 2.

#### 3.1 A Special Property of SRT-Feasible Task Sets

It follows from results in [16] that, under global scheduling,  $\tau$  is feasible (HRT or SRT) if and only if  $U \leq m$  holds and  $U_i \leq 1$  holds for each task  $\tau_i$ . That is, avoiding over-utilization is the key to ensuring feasibility. In Lemma 4 below, we show that the same is true for SRT-feasibility with arbitrary affinity masks, but the "no over-utilization" condition is more complicated. The next lemma, which is used in proving Lemma 4, refers to "uncompleted work": in a given schedule, the uncompleted work at time t is the total execution time of all jobs released prior to t minus the processing capacity already allocated to those jobs.

▶ Lemma 3. If the tardiness of every task in  $\tau$  is at most B in some schedule, then at any time instant in that schedule, the amount of uncompleted work is at most  $BU + 2\sum_i C_i$ .

**Proof.** If tardiness never exceeds B, then every job completes within B time units of its deadline, which for a job of task  $\tau_i$ , is within  $B+T_i$  time units of its release. Thus, at any time t, all jobs of task  $\tau_i$  released prior to time  $t-B-T_i$  are completed. During  $[t-B-T_i,t)$  task  $\tau_i$  may release at most  $\left\lceil \frac{B+T_i}{T_i} \right\rceil$  jobs. Thus, the amount of uncompleted work due to task  $\tau_i$  at time t is at most  $C_i \left\lceil \frac{B+T_i}{T_i} \right\rceil \leq C_i \left(2+\frac{B}{T_i}\right) = 2C_i + B\frac{C_i}{T_i} = 2C_i + BU_i$ . Summing over all tasks yields  $BU+2\sum_i C_i$ .

For  $\tau' \subseteq \tau$ , let  $U_{\tau'}$  denote  $\sum_{\tau_i \in \tau'} U_i$ . The "no over-utilization" condition we require is: **Utilization Balance:**  $\forall \tau' \subseteq \tau : U_{\tau'} \leq |\alpha_{\tau'}|$ .

▶ Lemma 4. If  $\tau$  is SRT-feasible, then it satisfies Utilization Balance.

**Proof.** Assume, contrary to the statement of the lemma, that a SRT-feasible task set  $\tau$  exists that violates Utilization Balance. Then, for some  $\tau' \subseteq \tau$ ,

$$|\alpha_{\tau'}| < U_{\tau'}. \tag{1}$$

Consider the following *periodic* release sequence for  $\tau$ : each task  $\tau_i$  in  $\tau$  releases jobs every  $T_i$  time units, starting at time 0, and each such job executes for  $C_i$  time units. Let S be the schedule with bounded tardiness for this release sequence mentioned in Lemma 3.

By the definition of  $\alpha_{\tau'}$ , all jobs of all tasks in  $\tau'$  are scheduled in S on cores from  $\alpha_{\tau'}$ . Let H be the hyperperiod of  $\tau$ . Then, for any integer k, the amount of work generated by  $\tau'$  over [0,kH) is  $\sum_{\tau_i \in \tau'} C_i \cdot \frac{kH}{T_i} = kHU_{\tau'} > kH|\alpha_{\tau'}|$ , where the last inequality follows from (1). Observe that  $kH|\alpha_{\tau'}|$  corresponds to the total available capacity over [0,kH) on cores in  $\alpha_{\tau'}$ . Thus, the uncompleted work at time kH in S is at least  $kH(U_{\tau'} - |\alpha_{\tau'}|)$ . This value grows unboundedly with increasing k, contradicting Lemma 3.

From results presented later, it will follow that Utilization Balance is a necessary and sufficient condition for SRT-feasibility (and also HRT-feasibility), *i.e.*, Lemma 4 can be strengthened by specifying "if and only if." When we henceforth wish to emphasize this usage of Utilization Balance, we will refer to it as the *UB Test*. Unfortunately, applying the UB Test by considering different subsets of tasks can require  $\Omega(2^n)$  time. However, the structure of this test is similar to the famous condition of Hall's Marriage Theorem [12], the proof of which involves examining maximal "edge matchings" in a graph. Such matchings can be determined by considering the Ford-Fulkerson max-flow algorithm and its correctness proof [11]. This connection to prior work (along with the existence of polynomial-time algorithms for max flow) motivates us to determine whether max flow can be used to efficiently determine SRT-feasibility.

#### 3.2 Max-Flow Feasibility Test

To cast checking feasibility as a max-flow problem, we define for any task set  $\tau$  a flow network  $FN(\tau)$  that is obtained from its affinity graph  $AG(\tau)$  via several steps. First, each edge  $(\tau_i, \pi_j)$  in  $AG(\tau)$  is viewed as a directed edge from  $\tau_i$  to  $\pi_j$  and given a capacity Z, where Z > m. Second, a source vertex s is added along with an edge  $(s, \tau_i)$  with capacity  $U_i$  for each vertex  $\tau_i$ . Finally, a sink vertex t is added along with an edge  $(\pi_j, t)$  with capacity 1 for each vertex  $\pi_j$ . Following conventional notation, we let f denote a flow that is defined with respect to  $FN(\tau)$ , with f(u, v) denoting the flow from vertex u to vertex v, and let |f| denote the value of the flow f (which equals the total flow coming from the source vertex s). (To avoid notational clutter, we do not parameterize f by  $\tau$ .)

- (a) Flow graph for the task set in Ex. 2. (b) A cut of the flow graph in inset (a).
- **Figure 3** Example flow graph and cut.
- ▶ **Example 5.** Fig. 3a shows the flow network corresponding to the affinity graph in Fig. 1b.
- ▶ **Lemma 6.** If Utilization Balance holds for  $\tau$  and f is a maximum flow, then |f| = U.

**Proof.** Assuming f is a maximum flow, by the Max-Flow/Min-Cut Theorem [11], |f| equals the capacity of a minimal cut. A *cut* is a partitioning of vertices that places the source and sink in different partitions. The *capacity* of a cut is simply the sum of the capacities of all edges that traverse the cut. Such an edge is called *crossing edge*. For example, Fig. 3b shows one of the many cuts that can be defined with respect to the flow network in Fig. 3a (the vertex sets  $V_V$ ,  $V_W$ , and  $V_X$  are discussed later). The capacity of this cut is  $U_4 + 2$ .

Let C be a cut with minimal capacity. If any edge of the form  $(\tau_i, \pi_j)$  is a crossing edge, then because its capacity exceeds m and the capacity of any edge is non-negative, the capacity of C exceeds m. This cannot be the case if C is minimal because the cut that places t and all other vertices in different partitions has capacity m. Thus, every crossing edge is of the form  $(s, \tau_i)$  or  $(\pi_j, t)$ . The cut show in Fig. 3b has this property, so the reader may wish to consult it for illustrative purposes hereafter.

Let  $V_W = \{\tau_i \mid (s,\tau_i) \text{ is a crossing edge}\}$ ,  $V_V = \{\tau_i \mid (s,\tau_i) \text{ is not a crossing edge}\}$ , and  $V_X = \{\pi_j \mid (\pi_j,t) \text{ is a crossing edge}\}$ . Then, the capacity of the cut C is  $\sum_{\tau_i \in V_W} U_i + |V_X|$ . As long as there are no crossing edges of the form  $(\tau_i,\pi_j)$ , all edges from tasks from  $V_V$  are to  $V_X$ . Thus,  $|V_X| \geq |\alpha_{V_V}|$ . Assuming Utilization Balance holds for  $\tau$ ,  $|V_X| \geq |\alpha_{V_V}| \geq \sum_{\tau_i \in V_V} U_i$ . Therefore, the capacity of C is at least  $\sum_{\tau_i \in V_W} U_i + |V_X| \geq \sum_{\tau_i \in V_W} U_i + \sum_{\tau_i \in V_V} U_i = \sum_{\tau_i \in \tau} U_i = U$ .

To summarize, any minimal cut has capacity at least U. However, the cut that places s and all other vertices in different partitions has capacity U, so the capacity of any minimal cut actually equals U. Thus, by the Max-Flow/Min-Cut Theorem, |f| = U.

It will follow from results presented next that showing that U is a maximum flow is an alternative way to test SRT-feasibility. We therefore refer to this alternative as the MF Test.

In fact, we are going to show that the UB Test and MF Test are each valid tests for *both* HRT- and SRT-feasibility. We do this by providing reasoning for the remaining links in the proof overview given earlier in Fig. 2. One of these links involves considering a HRT-feasibility test, denoted APA-Feas $(\tau, \pi)$ , presented by Baruah *et al.* [6]:

**APA-Feas** $(\tau, \pi)$  **Test:** Declare  $\tau$  to be HRT-feasible if and only if values exist for the variables  $x_{ij}$  satisfying:

$$\forall i: \sum_{j \in \alpha_i} x_{ij} = 1,$$
  $\forall j: \sum_i x_{ij} U_i \le 1, \text{ and }$   $\forall i, j: x_{ij} \ge 0.$ 

**Max-flow linear program.** The APA-Feas $(\tau, \pi)$  Test can be cast as a linear program (LP). Thus, to make a connection between it and our earlier results, we consider an LP-based implementation of the MF Test, which we refer to as the *LP-MF Test*:

Maximize 
$$\sum_{i} f(s, \tau_i)$$

Subject to:

Edge Constraints:

$$\begin{array}{lll} \text{C1:} & \forall i: \ 0 \leq f(s,\tau_i) \leq U_i & \{(s,\tau_i) \text{ edges}\} \\ \text{C2:} & \forall i,j: \ 0 \leq f(\tau_i,\pi_j) \leq Z & \{(\tau_i,\pi_j) \text{ edges}\} \\ \text{C3:} & \forall j: \ 0 \leq f(\pi_j,t) \leq 1 & \{(\pi_j,t) \text{ edges}\} \end{array}$$

Vertex Constraints:

C4: 
$$\forall i: f(s, \tau_i) = \sum_{j \in \alpha_i} f(\tau_i, \pi_j)$$
  $\{\tau_i \text{ vertices}\}$   
C5:  $\forall j: \sum_{i:(j \in \alpha_i)} f(\tau_i, \pi_j) = f(\pi_j, t)$   $\{\pi_j \text{ vertices}\}$ 

The edge constraints ensure that edge capacities are respected and the vertex constrains ensure that the flow into each non-source/sink vertex matches the flow out of that vertex.

We now show that if |f| = U holds, then the constraints in the LP above can be simplified, yielding constraints quite similar to those in the APA-Feas $(\tau, \pi)$  Test. In particular, if f is a maximum flow, then it still lies in the feasibility region of the simplified LP.

▶ **Lemma 7.** If f is a maximum flow and |f| = U, then the following conditions hold.

$$\forall i: \sum_{j \in \alpha_i} f(\tau_i, \pi_j) = U_i \tag{2}$$

$$\forall j: \sum_{i} f(\tau_i, \pi_j) \le 1 \tag{3}$$

$$\forall i, j, \text{ where } j \notin \alpha_i : f(\tau_i, \pi_j) = 0,$$
 (4)

assuming we assign  $f(\tau_i, \pi_j)$  to be 0 for  $j \notin \alpha_i$  (note that, if  $j \notin \alpha_i$ , then the edge  $(\tau_i, \pi_j)$  is not present in the flow network).

**Proof.** Let f be a maximum flow such that |f| = U. Because f is a maximum flow, it satisfies the constraints of the LP of the LP-MF Test. By Constraint C1,  $|f| = \sum_i f(s, \tau_i) \le \sum_i U_i = U$ . Because |f| = U, by the construction of the flow network,  $f(s, \tau_i) = U_i$  holds for each i. Thus, by Condition C4, (2) holds. Furthermore, by Conditions C3 and C5,  $\sum_{i:(j\in\alpha_i)} f(\tau_i, \pi_j) = f(\pi_j, t) \le 1$ , so (3) holds, assuming we assign  $f(\tau_i, \pi_j)$  to be 0 for  $j \notin \alpha_i$  as stated in the lemma. Note that such an assignment trivially satisfies (4).

#### 3.3 HRT- and SRT-Feasibility Equivalence

The following theorem summarizes the results above.

▶ Theorem 8. A task set  $\tau$  is SRT-feasible if and only if it is HRT-feasible. Moreover, the UB Test and the MF Test (and its LP counterpart, the LP-MF Test) are each both exact SRT- and HRT-feasibility tests.

**Proof.** By Lemma 4, if  $\tau$  is SRT-feasible, then it satisfies Utilization Balance, which by Lemma 6 implies that |f| = U holds, where f is a maximum flow. Thus, by Lemma 7, f satisfies Conditions (2)–(4). Now, defining,  $x_{ij}$  by  $x_{ij} = \frac{f(\tau_i, \pi_j)}{U_i}$ , Conditions (2)–(4) imply that all of the conditions of the APA-Feas $(\tau, \pi)$  Test are satisfied, so  $\tau$  is HRT-feasibile. As noted earlier, HRT-feasibility trivially implies SRT-feasibility. Thus, all links in the chain of reasoning depicted in Fig. 2 have been validated.

Remarks. Given Theorem 8, we will generally use the term "feasible" hereafter without qualifying whether we mean SRT- or HRT-feasibility.

Using a max flow algorithm from [14] with  $\tilde{O}(E\sqrt{(V)})$  time complexity, where V is the number of vertices and E is the number of edges in the flow network, the MF Test can be performed in  $\tilde{O}(mn\sqrt{m+n})$  time, since our flow network satisfies V=m+n+2 and  $E \leq mn$ . In contrast, the APA-Feas $(\tau, \pi)$  Test requires solving an LP with mn variables and n+m constraints, which requires  $\tilde{O}(mn(m+n)^{\omega+0.5})$  total time in the worst case [14], where  $2 < \omega < 2.4$  is the matrix multiplication constant [23]. Thus, the MF Test is considerably more efficient than the APA-Feas $(\tau, \pi)$  Test.

#### 4 **Loop-Free Affinities**

As a stepping stone towards defining algorithm AM-Red, we provide in this section an SRToptimal scheduler under the restriction that  $\alpha$  is loop-free. For any feasible task set  $\tau$  that this scheduler must correctly schedule, we fix f to be a maximum flow satisfying Conditions (2)-(4) of Lemma 7. Using this fixed f, the algorithm designed here seeks to ensure that each task  $\tau_i$  receives a long-term processor share on core  $\pi_j$  equal to  $f(\tau_i, \pi_j)$ .

**Share graph.** Note that  $f(\tau_i, \pi_j) = 0$  may hold even when  $\pi_j \in \alpha_i$ . In this case, even though task  $\tau_i$  is allowed to execute on core  $\pi_j$ , the share allocation defined above will preclude this from happening. Thus, while the affinity graph  $AG(\tau)$  includes the edge  $(\tau_i, \pi_j)$ , this edge really is not needed. To reflect this, we define a share graph  $SG(\tau)$  that is a subgraph of  $AG(\tau)$ . The two are the same except that, in  $SG(\tau)$ , any edge  $(\tau_i, \pi_i)$  in  $AG(\tau)$ for which  $f(\tau_i, \pi_i) = 0$  holds in the corresponding flow network  $FN(\tau)$  is removed. Note that  $SG(\tau)$  is loop-free if  $AG(\tau)$  is.

**Example 9.** Consider a task set  $\tau$  consisting of three tasks,  $\tau_1$ ,  $\tau_2$ , and  $\tau_3$ , scheduled on four cores,  $\pi_1$ ,  $\pi_2$ ,  $\pi_3$ , and  $\pi_4$ . If the max-flow values computed for  $\tau$  are as specified in Fig. 5a, then its share graph is as depicted in Fig. 5b. (Fig. 5 has several other insets that are considered later.)

#### 4.1 **Frames**

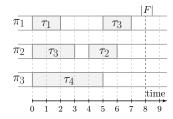
To realize the long-term processor shares that we want, we define allocations offline for a time interval called a frame. We denote the allocation function by F and the frame length by |F|. At runtime, we use F to perform allocations within each successive time window of length |F|. Formally, F is a mapping  $F:[0,|F|)\times\{\pi_1,...,\pi_m\}\to\{\emptyset,\tau_1,...,\tau_n\}$ . Informally, at each time instant within a widow of length |F|, F indicates which task is scheduled on each core (if core  $\pi_j$  is idle at time instant t, then  $F(t, \pi_j) = \emptyset$ ).

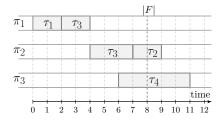
Let  $I_F(\tau_i, \pi_i)$  be the union of all maximal continuous intervals on core  $\pi_i$  allocated to task  $\tau_i$  by F. (Note that we use half-open intervals of the form [t, t').) Then, F is termed valid if the following conditions hold.

$$\forall i, j, j', j \neq j' : I_F(\tau_i, \pi_j) \cap I_F(\tau_i, \pi_{j'}) = \emptyset$$
(5)

$$\forall i, j : |I_F(\tau_i, \pi_j)| = f(\tau_i, \pi_j) \cdot |F| \tag{6}$$

<sup>&</sup>lt;sup>1</sup>  $\tilde{O}$  ignores logarithmic factors:  $\tilde{O}(g(n)) = O(g(n) \log^k g(n))$  for some natural number k.





- (a) A valid frame for  $\tau$  in Ex. 10.
- **(b)** A valid extended frame for  $\tau$  in Ex. 10.

**Figure 4** Example of a valid frame and a valid extended frame.

Informally, (5) implies that  $\tau_i$  cannot be allocated time on different cores simultaneously, while (6) implies that task  $\tau_i$  receives a total per-frame allocation of  $f(\tau_i, \pi_j) \cdot |F|$  on core  $\pi_j$ .

▶ Example 10. Consider the task set  $\tau$  from Ex. 2, with flow values  $f(\tau_1, \pi_1) = 1/4$ ,  $f(\tau_3, \pi_1) = 1/4$ ,  $f(\tau_2, \pi_2) = 1/4$ ,  $f(\tau_3, \pi_2) = 3/8$ , and  $f(\tau_4, \pi_3) = 5/8$ . A valid frame for  $\tau$  is depicted in Fig. 4a. Observe that |F| = 8 and (for example)  $\tau_3$ 's total allocation on core  $\pi_2$  is  $f(\tau_3, \pi_2) \cdot |F| = \frac{3}{8} \cdot 8 = 3$ .

**Extended frames.** To simplify the problem of defining a valid frame F for  $\tau$ , we introduce the concept of an extended frame E. E is a mapping similar to F but its length is not a priori bounded:  $E: [0, \infty) \times \{\pi_1, ..., \pi_m\} \to \{\emptyset, \tau_1, ..., \tau_n\}$ . We will show that a valid frame F can be obtained by first constructing E and then "wrapping" the allocations given by E to obtain F. First, we need to give validity conditions for E.

Let  $I_E(\tau_i, \pi_j)$  be the union of all maximal continuous intervals on core  $\pi_j$  allocated to task  $\tau_i$  by E. At the risk of a slight notational overload, let  $L_E(\tau_i) = \bigcup_j I_E(\tau_i, \pi_j)$  and  $L_E(\pi_j) = \bigcup_i I_E(\tau_i, \pi_j)$ . The conditions for E to be valid are as follows.

$$\forall i, j : I_E(\tau_i, \pi_j)$$
 is a single continuous interval (7)

$$\forall i: L_E(\tau_i) \text{ is a single continuous interval}$$
 (8)

$$\forall j: L_E(\pi_i) \text{ is a single continuous interval}$$
 (9)

$$\forall i, j : |I_E(\tau_i, \pi_j)| = f(\tau_i, \pi_j) \cdot |F| \tag{10}$$

$$\forall i, j_1, j_2 : I_E(\tau_i, \pi_{j_1}) \cap I_E(\tau_i, \pi_{j_1}) = \emptyset$$
(11)

▶ Example 11. Fig. 4b shows a valid extended frame for the task set in Ex. 10. Notice how each task is scheduled over a continuous interval of time. A task can migrate from one core to another during such an interval, but once it completes executing on a given core, it cannot execute on that core again. Also, once a core transitions from executing some task to being idle, it must stay idle.

Note that (10) ensures that the total allocation to  $\tau_i$  in E matches that in F. However, no constraints are placed on the length of E, so it may potentially contain many idle intervals.

Converting from extended frames to frames. The following lemma reduces the problem of defining a valid frame to that of defining a valid extended frame.

▶ **Lemma 12.** Assume that the extended frame E is valid. Furthermore, if  $E(t, \pi_j) \neq \emptyset$  for some t, then define  $F(t \text{ mod } |F|, \pi_j) = E(t, \pi_j)$ . Then, F is valid.

**Proof.** First, note that F is well-defined: by (9) and (10), allocations on  $\pi_j$  obtained from E occur within an interval of length  $\sum_i f(\tau_i, \pi_j) \cdot |F| \leq |F|$ , where the latter inequality follows

#### Algorithm 1 Extended-Frame Builder

```
Require: cores order (\prec), tasks order for each for \pi_j (\stackrel{\pi_j}{\longrightarrow})

Ensure: extended frame

1: for \pi_j \in \text{cores order do}

2: for \tau_i \in \text{tasks order for } \pi_j \text{ do}

3: if \tau_i is the first task allocated on \pi_j then

4: start \leftarrow \text{end of the last allocation interval for } \tau_i \text{ if one exists, else 0}

5: else

6: start \leftarrow \text{end of the last allocation interval on core } \pi_j

7: define the allocation interval for \tau_i on \pi_j to be [start, start + f(\tau_i, \pi_j))
```

from (3). Thus, if  $\pi_j$  is allocated (not idle) at distinct time instants t and t' by E, then  $t \mod |F| \neq t' \mod |F|$ .

To show that F is valid, we must show that (5) and (6) hold. By (7), (8), and (10),  $L_E(\tau_i)$  has length  $(\sum_{j \in \alpha_i} f(\tau_i, \pi_j)) \cdot |F| = U_i \cdot |F| \le |F|$ , where the first equality follows from (2). Thus, a similar argument as given in the first paragraph of the proof can be applied to show that (5) holds. As for (6), it follows directly from (10).

Constructing a valid extended frame. Given Lemma 12, we can focus our attention now on constructing a valid extended frame. To motivate some of the issues that arise in doing so, we first consider an example.

**Example 13.** We can compute a valid extended frame E for the task set in Ex. 9 as follows. Consider the cores in order, and for each core  $\pi_j$ , consider the tasks connected by an edge to  $\pi_j$  in turn. This ordering is illustrated in Fig. 5c as an "outer" ordering of cores and an "inner" ordering of tasks. Given this ordering, we can apply the simple scheme in Alg.1 to obtain E. First, consider core  $\pi_1$  and tasks  $\tau_1$  and  $\tau_2$  in turn. Allocate them shares of 1/3 and 2/3, respectively, on  $\pi_1$ , starting from time 0. Now, move on to core  $\pi_2$  and consider the tasks  $\tau_1$  and  $\tau_3$  in turn. Allocate them shares of 1/3 and 2/3, respectively, on  $\pi_2$ , but this time starting from the end of  $\tau_1$ 's allocation on  $\pi_1$  (so that (11) is not violated). Continue to consider cores and tasks in this manner until all tasks have received their needed share allocations. The resulting extended frame E that is constructed in shown in Fig. 5d.

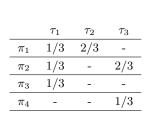
Ordering cores and tasks. The determination of E worked out easily in Ex. 13 because we conveniently ordered cores and tasks in way to make this happen. Other orderings could be problematic. For example, with the core ordering  $\pi_1, \pi_4, \pi_2, \pi_3$ , the obtained extended frame E would not be valid because (8) would be violated for  $\tau_3$ , as illustrated in Fig. 5e. Properly ordering cores is not enough. For example, had we kept the original core ordering but changed the ordering of tasks on core  $\pi_2$  to  $\tau_3, \tau_1$ , then the obtained extended frame would again violate (8), this time for  $\tau_1$ , as illustrated in Fig. 5f. These examples show that properly ordering cores and tasks is crucial for Alg.1 to be correct.

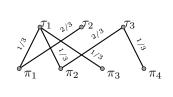
**Proper orderings.** To correctly apply Alg.1 in a general way, we define a total order  $\prec$  on cores, and for each core  $\pi_j$ , a total order  $\xrightarrow{\pi_j}$  on a certain subset of tasks  $\tau^{\pi_j}$ . We say that a task  $\tau_i$  is a non-first task on core  $\pi_j$  if and only if there exists  $\tau_{i'}$ , where  $i' \neq i$ , such that  $\tau_{i'} \xrightarrow{\pi_j} \tau_i$ . These orders are called proper if the following conditions hold.

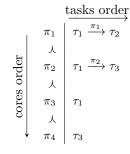
$$\forall i, j: \ \tau_i \in \tau^{\pi_j} \text{ if and only if } f(\tau_i, \pi_j) > 0$$
 (12)

$$\forall i: \text{ task } \tau_i \text{ can be non-first task on at most one core}$$
 (13)

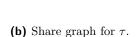
$$\forall i: \text{ if task } \tau_i \text{ is non-first on core } \pi_j, \text{ then for every } \pi_{j'} \prec \pi_j, \tau_i \notin \tau^{\pi_{j'}}$$
 (14)



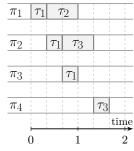


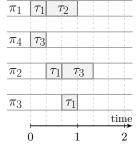


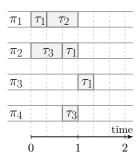
(a) Max-flow values for  $\tau$ .



(c) Ordering of cores and tasks.







- (d) Valid extended frame.
- (e) Invalid extended frame (improper core ordering).
- (f) Invalid extended frame (improper task ordering).
- **Figure 5** Illustrations pertaining to task set  $\tau$  from Exs. 9 and 13.
- ▶ **Example 14.** It is straightforward to verify that (12)–(14) hold for the orders in Fig. 5c. Thus, these orders are proper.
- ▶ Lemma 15. If Alg. 1 is provided proper orders, then a valid extended frame is returned.

**Proof.** Alg. 1 places at most one allocation interval for each task on any core, so (7) holds. The algorithm also ensures that if multiple allocation intervals are placed on some core, each successive interval begins when the immediately prior one ends, so (9) holds. By (12), the algorithm places an interval of size  $f(\tau_i, \pi_j)$  on  $\pi_j$  whenever  $f(\tau_i, \pi_j) > 0$  holds, so (10) holds. (13) and (14) ensure that, if task  $\tau_i$  has allocation intervals placed on different cores, it can be a non-first task only on the first (by  $\prec$ ) of these cores. Thus, Alg. 1 ensures that these allocation intervals are contiguous, so (8) and (11) hold.

Generating proper orders. The final issue that remains is actually generating proper cores/tasks orders. For this, we provide Alg. 2. (Note that, if  $SG(\tau)$  is not connected, then we assume that the BFS routine searches every connected component.)

▶ **Theorem 16.** If  $SG(\tau)$  is loop-free, then Alg. 2 produces orders that are proper.

**Proof.** Line 4 of Alg. 2 ensures that (12) holds. In the rest of the proof, we verify the remaining properies, (13) and (14). We assume that all vertices and edges referenced in verifying these properties are part of the same connected component of  $SG(\tau)$ .

Assume, to the contrary of (13), that  $\tau_i$  is a non-first task on two distinct cores,  $\pi_j$  and  $\pi_{j'}$ . Then,  $SG(\tau)$  has edges  $(\tau_i, \pi_j)$  and  $(\tau_i, \pi_{j'})$ . Furthermore,  $\pi_j$  and  $\pi_{j'}$  were discovered before  $\tau_i$ . Without loss of generality, assume that  $\pi_j$  was discovered first. Then, there exists

# Algorithm 2 Proper-Orders Generator

```
Require: loop-free share graph SG(\tau)

Ensure: proper core/task orders

1: run BFS(SG(\tau)) \qquad \qquad \triangleright breadth-first search of SG(\tau)

\qquad \qquad \triangleright Every connected component of SG(\tau) is searched starting with an arbitrary vertex

2: define cores order, \prec, by ordering cores according to their BFS discovery times

3: for \pi_j \in cores order do

4: \tau^{\pi_j} \leftarrow \{\tau_i \mid f(\tau_i, \pi_j) \neq 0\}

5: define the task order for \pi_j, \xrightarrow{\pi_j}:

6: if \pi_j is discovered in BSF by traversing the edge (\tau_i, \pi_j) then place \tau_i first

7: order other tasks in \tau^{\pi_j} arbitrarily after \tau_i (if it exists)
```

#### Algorithm 3 Loop-Free Scheduler

```
Require: \tau, frame\_len
 1: function FeasibilityCheck(\tau)
                                                                                    ▷ described in Sec. 3.2
        construct flow network FN(\tau)
 2:
        compute max flow f with respect to FN(\tau)
 3:
        if |f| = U then
 4:
                                                                                                 \triangleright MF Test
 5:
            return f
 6:
        else
                                                                                            \,\vartriangleright\,\tau is infeasible
 7:
            return \perp
 8: function GenerateFrame(\tau, frame_len, f)
                                                                                    \triangleright described in Sec. 4.1
       run Alg. 2 to get proper core/task orders
        run Alg. 1 to build a valid extended frame E
10:
        apply the transformation of Lemma 12 to obtain a valid frame F with |F| = frame len
11:
        return F
13: function Scheduler(\tau, frame_len)
        f \leftarrow \text{FeasibilityCheck}(\tau)
14:
        if f \neq \bot then
15:
            F \leftarrow \text{GENERATEFRAME}(\tau, frame\_len, f)
16:
            repeat the allocations in F every |F| time units, letting the jobs of each task \tau
17:
                execute within the allocation intervals for \tau in release-time order
```

a path  $\pi_j \rightsquigarrow \pi_{j'}$  that does not include  $\tau_i$ . Thus, we have cycle,  $\pi_j \rightsquigarrow \pi_{j'} \to \tau_i \to \pi_j$ , which is a contradiction.

Finally, assume to the contrary of (14), that  $\tau_i$  is non-first on core  $\pi_j$ , but core  $\pi_{j'}$  exists such that  $\pi_{j'} \prec \pi_j$  and  $\tau_i \in \tau^{\pi_{j'}}$ . Because  $\pi_{j'} \prec \pi_j$  holds,  $\pi_{j'}$  was discovered before  $\pi_j$ . Because  $\tau_i \in \tau^{\pi_{j'}}$ , the edge  $(\tau_i, \pi_{j'})$  exists. Cores are not connected by edges in  $SG(\tau)$ , so these facts imply that  $\tau_i$  was discovered before  $\pi_j$ . Because  $\tau_i$  was selected as a non-first task on core  $\pi_j$ , the edge  $(\tau_i, \pi_j)$  exists. It follows that  $\pi_j$  would have been discovered by traversing that edge, making  $\tau_i$  the first-ordered task on core  $\pi_j$ , which is a contraction.

#### 4.2 Scheduler

We summarize our results so far by presenting Alg. 3, our algorithm for scheduling task sets with loop-free affinity masks. We present analysis pertaining to this scheduler below, after first providing an example that illustrates how it works.

▶ Example 17. Consider the task set from Ex. 2 with  $f(\tau_i, \pi_j)$  values from Ex. 10 and the (valid) frame F shown in Fig. 4a. Fig. 6 shows how several jobs of the sporadic migrating task  $\tau_3$  are scheduled under Alg. 3 assuming  $T_3 = 1.25|F|$ . The core label within each allocation interval indicates the core upon which  $\tau_3$  is scheduled.

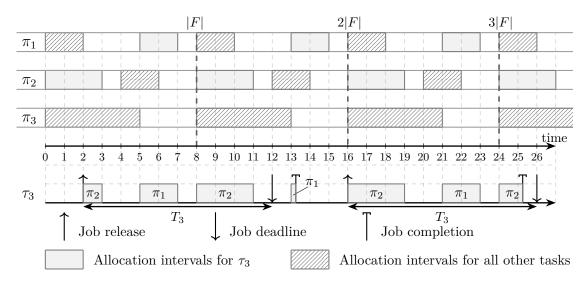


Figure 6 Example schedule under Alg. 3.

Frame-based schedulers were first studied years ago [19]. More recently, several frame-based semi-partitioned schedulers have been proposed, [5, 13, 21, 24], but none support affinities. Our frame-based scheduler draws inspiration from one of these [24], but that scheduler is directed at heterogeneous multiprocessors.

#### 4.3 Analysis and Optimality

In this section, we analyze Alg. 3 from the perspectives of task migrations (both the number of migrating tasks and migration frequency), time complexity, and optimality.

**Migrations.** We call a task *migrating* if it has allocations on multiple cores and *fixed* otherwise. Let M denote the number of migrating tasks in Alg. 3. The following theorem shows that Alg. 3 limits M in accordance with the idea of semi-partitioned scheduling, where the goal usually is to have M = O(m).

#### ▶ Theorem 18. $M \leq m-1$ .

**Proof.** By assumption,  $SG(\tau_i)$  is loop-free, which implies that it has at most n+m-1 edges. M migrating tasks have at least 2M incident edges. Thus, the number of edges incident upon fixed tasks is at most n+m-1-2M. Each fixed task has at least one incident edge (if the task set is feasible). It follows that at most n+m-1-2M tasks are fixed. Because n is the total number of tasks, we therefore have  $n \leq M+n+m-1-2M$ , implying  $M \leq m-1$ .

We now prove several migration-related bounds, all of which are *tight*, *i.e.*, task sets can be defined for which exactly these bounds hold. In proving these bounds, we let F be the (valid) frame used by Alg. 3, and let  $\deg(\tau_i)$  denote the degree of  $\tau_i$  in  $SG(\tau)$ . When we refer below to an *allocation interval* for a task  $\tau_i$ , we mean a maximal continuous interval during which  $\tau_i$  is allocated capacity on one core.

#### ▶ **Lemma 19.** Task $\tau_i$ has at most $\deg(\tau_i) + 1$ allocation intervals in F.

**Proof.** It is straightforward to show that, in the valid extended frame E that is used to obtain F,  $\tau_i$  has at most  $\deg(\tau_i)$  allocation intervals. Using (7)–(11), it is easy to show that these intervals occupy a continues time window of length at most |F|. Thus, when this

interval is "wrapped" (see Lemma 12) to produce F, at most one of these intervals is split into two subintervals. The stated bound follows.

▶ **Theorem 20.** The overall number of migrations within F is at most 2m-2.

**Proof.** Let  $\tau^M$  denote the set of migrating tasks. To compute the overall number of migrations, we consider only these tasks. By Lemma 19, each such task has at most  $\deg(\tau_i)+1$  allocation intervals in F. This yields a bound of  $\deg(\tau_i)$  for the number of migration by  $\tau_i$  in F because a task's first allocation interval does not entail a migration. Recall (from the proof of Theorem 18) that  $SG(\tau)$  has at most n+m-1 edges. Thus, the total number of migrations within F is at most  $\sum_{\tau_i \in \tau^M} \deg(\tau_i) = \text{no.}$  of edges in  $SG(\tau) - \sum_{\tau_i \notin \tau^M} \deg(\tau_i) \leq n+m-1-(n-M) = m-1+M$ , which by Theorem 18, is at most 2m-2.

▶ **Theorem 21.** Within any continuous time interval of length L, task  $\tau_i$  has at most  $\lceil L/|F| \rceil \cdot \deg(\tau_i)$  migrations, and the overall number of migrations is at most  $\lceil L/|F| \rceil \cdot (2m-2)$ .

**Proof.** As discussed in the proof of Lemma 19, when "wrapping" the extended frame E to get F, a task  $\tau_i$  has at most one allocation interval that is split. If it has zero, then its first and last allocations per frame are for different cores, while if it has one, they are for the same core. Thus, with zero split allocations, inter-frame migrations can occur, while with one, they cannot. Hence, reasoning as in the proofs of Lemma 19 and Theorem 20, when accounting for both intra- and inter-frame migrations, we have at most  $\deg(\tau_i)$  migrations for  $\tau_i$  per frame and at most 2m-2 per frame overall. Thus, within L,  $\tau_i$  experiences at most  $\lceil L/|F| \rceil \cdot \deg(\tau_i)$  migrations, and the overall number of migrations is at most  $\lceil L/|F| \rceil (2m-2)$ .

Time complexity. It is straightforward to show that Algs. 1 and 2 each can be implemented in O(m+n) time. In contrast, the most efficient known max-flow algorithms require superlinear time. Thus, the time complexity required by Alg. 3 to find a valid frame is dominated by that of the max-flow algorithm that is used. Lee *et al.* [14] have presented a max-flow algorithm that has time complexity  $\tilde{O}(mn\sqrt{m+n})$  for an arbitrary  $AG(\tau)$  and  $\tilde{O}((m+n)^{3/2})$  in our setting (the number of edges in a loop-free  $AG(\tau)$  is limited by (m+n-1)), giving us the following.

▶ **Theorem 22.** For any feasible task set  $\tau$ , Alg. 3 can produce a valid frame F in  $\tilde{O}(mn\sqrt{m+n})$ ; if  $\alpha_{\tau}$  is loop-free, it requires  $\tilde{O}((m+n)^{3/2})$  time.

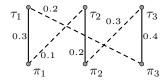
Bonifaci et al. [8] claim that semi-partitioned and hierarchical scheduling with affinity masks is NP-hard to approximate. Their work does not contradict Theorem 22 because it is directed at a completely different context, namely, one-shot jobs and makespan minimization.

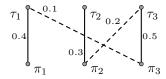
**Optimality.** The next lemma shows that each task receives a long-term processor share equal to its utilization. We use this lemma in showing that Alg. 3 is optimal below.

▶ Lemma 23. Task  $\tau_i$  receives a total allocation of at least  $U_i \cdot |F| \cdot \lfloor L/|F| \rfloor$  during any interval of length L.

**Proof.** During an interval of length |F|,  $\tau_i$  receives an allocation of  $U_i \cdot |F|$ . (Note that, if such an interval begins within a frame, then the missing part of that frame is compensated for by the beginning of the next frame—recall that all frames are identical). During an interval of length L, at least |L/|F| complete intervals of length |F| occur.

▶ **Theorem 24.** Alg. 3 is SRT-optimal and ensures a tardiness bound of |F|. Moreover, if |F| divides all periods, then it is HRT-optimal.





- (a) Share graph before loop removal.
- (b) Share graph after loop removal.

Figure 7 Ex. 25. (For clarity, task-to-core edges are solid and core-to-task edges are dashed.)

**Proof.** Consider a job  $J_{i,s}$  that is released at time  $t_r$  and has a deadline at time  $t_d$ . Let t be the last time instant at or before  $t_r$  such that no job of  $\tau_i$  released prior t is unfinished at t. Let k be number of jobs of  $\tau_i$  released in  $[t, t_d)$ . Then, because  $\tau_i$  has a deadline at  $t_d$ , we have  $t_d - t \ge kT_i$ . Thus, by Lemma 23, the total processor allocation received by  $\tau_i$  during  $[t, t_d + |F|)$  is at least

$$\begin{split} U_i \cdot |F| \cdot \left\lfloor \frac{t_d + |F| - t}{|F|} \right\rfloor &\geq U_i \cdot |F| \cdot \left\lfloor \frac{k \cdot T_i}{|F|} + 1 \right\rfloor \geq U_i \cdot |F| \cdot \left\lceil \frac{k \cdot T_i}{|F|} \right\rceil \\ &\geq U_i \cdot |F| \cdot \left( \frac{k \cdot T_i}{|F|} \right) = k \cdot C_i. \end{split}$$

This implies that  $J_{i,s}$  completes by time  $t_d + |F|$ , *i.e.*, Alg. 3 is SRT-optimal and ensures a tardiness bound of |F|.

If |F| divides all task periods, then similar reasoning can be applied, but this time with respect to the interval  $[t, t_d)$ . Because |F| divides  $T_i$ , by Lemma 23, the total processor allocation received by  $\tau_i$  during this interval is at least

$$\begin{aligned} U_i \cdot |F| \cdot \left\lfloor \frac{t_d - t}{|F|} \right\rfloor &\geq U_i \cdot |F| \cdot \left\lfloor \frac{k \cdot T_i}{|F|} \right\rfloor = U_i \cdot |F| \cdot \frac{k \cdot T_i}{|F|} \\ &= U_i \cdot |F| \cdot \left( \frac{k \cdot T_i}{|F|} \right) = k \cdot C_i. \end{aligned}$$

This implies that  $J_{i,s}$  completes with zero tardiness and that Alg. 3 is HRT-optimal.

#### 5 Arbitrary Affinity Masks

In this section, we show that arbitrary affinities can be dealt with by eliminating any loops that may exist in  $SG(\tau)$ . Because  $SG(\tau)$  is bipartite, any such loop is of the form  $\tau_{i_1} \to \pi_{j_1} \to \tau_{i_2} \to \pi_{j_2} \to \dots \to \pi_{j_k} \to \tau_{i_1}$ . Let  $f_m = \min(f(\tau_{i_1}, \pi_{j_2}), f(\tau_{i_2}, \pi_{j_1}), \dots, f(\tau_{i_1}, \pi_{j_k}))$ , i.e.,  $f_m$  is the minimal f value of any task-to-core edge in this loop. Then, we can eliminate the loop by decreasing the f value of each core-to-task edge by  $f_m$  and by increasing the f value of each task-to-core edge by  $f_m$ . This eliminates all loop edges with f values of  $f_m$ .

▶ **Example 25.** Applying this loop-removal procedure to the share graph in Fig. 7a, we have  $f_m = 0.1$ , and the share graph in Fig. 7b results.

Now that we have a procedure for eliminating loops, we need a means for finding them. A breadth-first-search routine can do this, as seen in Alg. 4. It is easy to see that this algorithm produces a loop-free share graph: it does not add any new edges, so it cannot create any new loops; also, any loop initially in  $SG(\tau)$  will be found by BFS and removed.

#### Algorithm 4 Affinity-reduction algorithm

```
Require: f

1: construct SG(\tau)

2: for \tau_i \in \tau do

3: while true do

4: Run BFS(\tau_i)

5: if BFS found a loop in SG(\tau) then

6: remove edge from the loop using the procedure illustrated in Ex. 25

7: else break while \Rightarrow no more loops in this connected component of SG(\tau)

8: return obtained f
```

#### Algorithm 5 AM-Red Scheduler

```
Require: \tau, frame_len
 1: function FeasibilityCheck(\tau)
                                                                                      \triangleright described in Sec. 3.2
 2:
        construct flow network FN(\tau)
 3:
        compute max flow f with respect to FN(\tau)
        if |f| = U then return f
                                                                                                   \triangleright MF Test
 4:
                                                                                             \triangleright \tau is infeasible
 5:
        else return \perp
 6: function GenerateFrame(\tau, frame\_len, f)
                                                                                      \triangleright described in Sec. 4.1
 7:
        run Alg. 2 to get proper core/task orders
 8:
        run Alg. 1 to build a valid extended frame {\cal E}
 9:
        apply the transformation of Lemma 12 to obtain a valid frame F with |F| = frame\_len
10:
        return F
11: function GetFlow(\tau)
        f \leftarrow \text{FeasibilityCheck}(\tau)
        run Alg. 4 to obtain loop-free f values
13:
14:
        return f
15: function Scheduler(\tau, frame\_len)
16:
        f \leftarrow \text{GetFlow}(\tau)
17:
        if f \neq \bot then
18:
            F \leftarrow \text{GENERATEFRAME}(\tau, frame\_len, f)
            repeat the allocations in F every |F| time units, letting the jobs of each task \tau
19:
                execute within the allocation intervals for \tau in release-time order
```

To determine the time complexity of Alg. 4, note that each invocation of the BFS routine requires O(E) time, where E is the number of edges in the initial graph  $SG(\tau)$ . In our context, E is upper bounded by the number of edges in  $AG(\tau)$ , which is O(mn). O(E+n) BFS-routine invocations occur, because each invocation removes at least one edge or moves to a new task. The edge-removal procedure itself (illustrated in Ex. 25) requires O(m) time. From this discussion, we have the following.

▶ **Theorem 26.** Alg. 4 transforms  $SG(\tau)$  into a loop-free graph. Its time complexity is  $O(m^2n^2)$  generally, and  $O(n^2)$  if the number of edges in  $AG(\tau)$  is linear.

**Algorithm AM-Red.** We are finally in a position to present the main contribution of this paper, algorithm AM-Red. As seen in Alg. 5, it is obtained by applying the various algorithms presented in this paper as building blocks in the expected way. The following theorem follows from the properties of the algorithms employed (particularly, properties stated as theorems).

▶ Theorem 27. For any feasible task set  $\tau$ , AM-Red produces a valid frame F in  $O(m^2n^2)$  time; it ensures at most m-1 tasks migrate and at most 2m-2 migrations occur per frame; it is SRT-optimal with a tardiness bound of |F|; if F divides all periods, it is also HRT-optimal.

#### 6 Hierarchical Masks

In this section, we show that the relatively high time complexity of the MF Test and affinity reduction (Alg. 4), which dominate the time complexity of AM-Red (Alg. 5), can be avoided if  $\alpha_{\tau}$  is hierarchical. To facilitate showing this, we assume that tasks are indexed such that  $i \leq j \Rightarrow |\alpha_i| \leq |\alpha_j|$ . We call this ordering canonical order. For now, we assume this ordering is initially provided. Later, we consider the cost of establishing it if not initially provided.

We now establish a simpler feasibility test when  $\alpha_{\tau}$  is hierarchical. To facilitate our description of this test, we introduce some new terminology. We say that task  $\tau_i$  is nested within task  $\tau_j$  if and only if  $i \leq j$  and  $\alpha_i \subseteq \alpha_j$ . It is easy to see that the "nested within" relation is transitive. We denote the set of tasks nested within  $\tau_i$  as  $N_i$  (note that  $\tau_i \in N_i$ ). Observe that, if  $\tau_i$  is nested within  $\tau_j$ , then  $N_i \subset N_j$  by transitivity. For any task  $\tau_i$ , we define its utilization closure as  $U_i^* = \sum_{\tau_j \in N_i} U_j$ . We call a task  $\tau_i$  maximal if it is not nested within any other task  $\tau_j$  with the same affinity mask, where  $j \geq i$ . For any task  $\tau_i$ , we let  $A_i$  denote the set of tasks with the same affinity mask (note that  $\tau_i \in A_i$ ); we say that these tasks agree with  $\tau_i$ . Note that the last task in  $A_i$  (in canonical order) is maximal. For  $\tau' \subseteq \tau$ , we let  $X(\tau')$  denote the set of all maximal tasks in  $\bigcup_{\tau_i \in \tau'} A_i$ , and we let  $\hat{X}(\tau')$  denote those tasks in  $X(\tau)$  at the "top" of the nesting hierarchy, i.e.,  $\hat{X}(\tau') = \{\tau_i : \tau_i \in X(\tau) \land \tau_i$  is not nested within any task in  $X(\tau)$ . Note that distinct tasks in  $\hat{X}(\tau')$  have disjoint masks.

- ▶ Example 28. Consider task set  $\tau$  from Ex. 2. Its affinity graph is shown in Fig. 1b. Consider the canonical order  $\tau_1, \tau_4, \tau_2, \tau_3$ . (When reasoning abstractly, we assume this ordering is consistent with task indices, as noted above.) Then,  $N_1 = \{\tau_1\}$ ,  $N_2 = \{\tau_1, \tau_2\}$ ,  $N_3 = \{\tau_1, \tau_2, \tau_3\}$ , and  $N_4 = \{\tau_4\}$ . Also,  $A_1 = \{\tau_1\}$ ,  $A_2 = A_3 = \{\tau_2, \tau_3\}$ , and  $A(\tau_4) = \{\tau_4\}$ . Tasks  $\tau_1, \tau_3, \tau_4$  are maximal, but task  $\tau_2$  is not since it is nested within  $\tau_3$ . For  $\tau' = \{\tau_1, \tau_2, \tau_4\}$ ,  $X(\tau') = \{\tau_1, \tau_3, \tau_4\}$  and  $\hat{X}(\tau') = \{\tau_3, \tau_4\}$ . Note that  $\tau_3$  and  $\tau_4$  have disjoint masks.
- ▶ **Lemma 29.** For any  $\tau_i \in \tau'$ ,  $\tau_i$  is nested within some task in  $X(\tau')$  and also  $\hat{X}(\tau')$ .

**Proof.** Any task in  $X(\tau)$  is nested within some task in  $\hat{X}(\tau')$ , so we can limit attention to  $X(\tau)$ . If  $\tau_i \in \tau'$  and  $\tau_i$  itself is not maximal, then it is nested within another task  $\tau_j$  with the same mask that is maximal. Because  $\tau_i \in \tau'$ ,  $\tau_j \in X(\tau')$ .

The simplified feasibility condition we require is as follows.

**Nested Balance:** For any maximal task  $\tau_i : U_i^* \leq |\alpha_i|$ .

▶ Lemma 30. In  $\alpha_{\tau}$  is hierarchical, then Utilization Balance (UB) and Nested Balance (NB) are equivalent.

**Proof.** It is straightforward to show UB  $\Rightarrow$  NB: if UB holds, then by considering  $\tau' = N_i$  in that condition, NB easily follows. In the rest of the proof, we focus on showing NB  $\Rightarrow$  UB. Consider any  $\tau' \subseteq \tau$  and any task  $\tau_i \in \tau'$ . By Lemma 29,  $\tau_i$  is nested within some task in  $\hat{X}(\tau')$ . Thus,  $\tau' \subseteq \bigcup_{\tau_j \in \hat{X}(\tau')} N_j$ , which implies  $U_{\tau'} \leq \sum_{\tau_j \in \hat{X}(\tau')} U_j^*$ . By NB,  $\sum_{\tau_j \in \hat{X}(\tau')} U_j^* \leq \sum_{\tau_j \in \hat{X}(\tau')} |\alpha_j|$ . As observed earlier, all tasks from  $\hat{X}(\tau')$  have disjoint masks. Moreover, the cores included in these masks are exactly the same as those included in  $\alpha_{\tau'}$ . Hence,  $\sum_{\tau_j \in \hat{X}(\tau')} |\alpha_j| = |\alpha_{\tau'}|$ . Putting these facts together, we have  $U_{\tau'} \leq \sum_{\tau_j \in \hat{X}(\tau')} U_j^* \leq |\alpha_{\tau'}|$ .

#### Algorithm 6 Nested Balance Test

```
Require: \tau (in canonical order)
 1: cap[j] = 0
                                                                    > capacity used on each core, bounded by one
 2: for \tau_i \in \tau do
                                                                       \triangleright r is remaining unallocated utilization of \tau_i
 3:
         r \leftarrow U_i
         while \pi_i from \alpha_i with cap[j] < 1 exists and r > 0 do
 4:
 5:
              f(\tau_i, \pi_j) \leftarrow \min(1 - cap[j], r)
 6:
              r \leftarrow r - f(\tau_i, \pi_j)
                                                                 \triangleright we used all available utilization of \tau_i (if r=0),
 7:
              cap[j] \leftarrow cap[j] + f(\tau_i, \pi_j)
                                                                            \triangleright or we filled core \pi_j fully (if cap[j] = 1)
         if r > 0 then return \perp
                                                                                            ▷ Nested Balance is violated
 8:
 9:
         else any non-defined f(\tau_i, \pi_j) value is considered to be 0
```

Having established a simpler feasibility condition, it remains to show it can be efficiently computed. Alg. 6 does this while also returning all needed non-zero  $f(\tau_i, \pi_j)$  values.

**Analysis.** We now show that Alg. 6 is correct and analyze its time complexity.

▶ **Theorem 31.** Alg. 6 returns  $f(\tau_i, \pi_j)$  values satisfying (2)-(4) iff  $\tau$  satisfies Nested Balance.

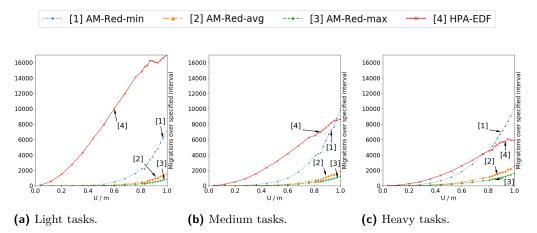
**Proof.** Establishing the algorithm's correctness when it does not return  $\bot$  for  $\tau$  is straightforward, so we focus on the other possibility, *i.e.*, it returns  $\bot$  when considering some task  $\tau_i$  in  $\tau$ . Because tasks are processed in canonical order, by the time  $\tau_i$  is considered, all tasks from  $N_i/\{\tau_i\}$  have already been dealt with and no task with a larger mask has yet been considered. Thus, the cores in  $\alpha_i$  could only have been allocated to tasks in  $N_i$ . If we cannot allocate  $\tau_i$ , then  $U_i^* = \sum_{\tau_j \in N_i} U_j > |\alpha_i|$ . If  $\tau_i$  is maximal, then Nested Balance is violated (and hence by Theorem 8 and Lemma 30,  $\tau$  is infeasible). Otherwise, there exists a maximal task  $\tau_k$  with the same mask as  $\tau_i$  but ordered after  $\tau_i$ . In this case,  $N_i \subset N_k$ , so  $U_k^* > U_i^* > |\alpha_i| = |\alpha_k|$ . Therefore, Nested Balance is violated in this case as well.

▶ Theorem 32. Alg. 6 completes in O(m+n) time. Thus, if  $\alpha_{\tau}$  is hierarchical, tasks are indexed in canonical order, and Alg. 6 is used in place of GETFLOW function in AM-Red, then AM-Red produces a valid frame in O(m+n) time.

**Proof.** During each **while**-loop iteration, either some core  $\pi_j$  becomes fully allocated (cap[j] becomes one), or the current task becomes fully allocated (r becomes zero), or  $\bot$  is returned. Each of these possibilities may happen only once. This implies that the total number of iterations of Alg. 6 is at most m + n.

If  $\alpha_{\tau}$  is hierarchical, then only O(m) unique affinity masks may exist [20, Theorem 3.5]. Thus, only O(m+n) space is required to provide canonically ordered tasks as input, as each task merely requires a pointer to one of the O(m) masks that may exist. If canonical order cannot be pre-assumed, then tasks must be sorted so they are so ordered. This can be done in  $O(m \log m + n)$  time: the masks themselves can be sorted in  $O(m \log m)$  time, and all per-task mask pointers can be updated in O(n) time. If one takes the core count m to be a constant (which is a very reasonable assumption), then the masks can be sorted in O(1) time and the total time complexity required to put tasks into canonical order is only O(n).

Example task sets can easily be constructed that have  $\Omega(m)$  distinct hierarchical masks. Because any feasibility test must consider these masks and all tasks, it follows that O(m+n) time complexity for testing feasibility is asymptotically optimal.



**Figure 8** Exp. 1 (hierarchical masks): total number of migrations under AM-Red and HPA-EDF (assuming periodic releases), averaged over the generated task sets, as a function of relative system utilization.

## 7 Experimental Evaluation

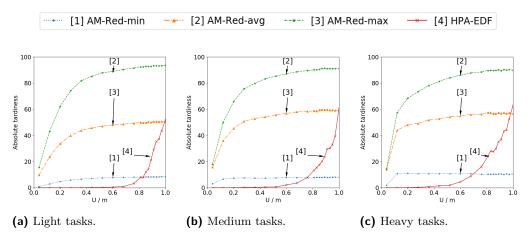
Given the limited range of prior related work, we compared AM-Red to HPA-EDF, a scheduler proposed in [7]. In this comparison, we were forced to limit attention to hierarchical masks, as HPA-EDF requires this. It is worth noting that AM-Red has linear offline time complexity and requires O(1) time per scheduling decision, while HPA-EDF requires O(m) time per job release and  $O(\log n + m^2)$  time per job completion. Thus, HPA-EDF has higher runtime overhead and this could negatively impact tardiness in practice. However, our comparison is based on computed schedules, and not actual scheduler implementations, so such runtime effects are not captured. As such, our comparison actually affords HPA-EDF an advantage.

Because |F| is tunable parameter, we considered three different ways of choosing it, which resulted in three variants of AM-Red: AM-Red-min, for which  $|F| = \min\{T_i\}$ ; AM-Red-avg, for which  $|F| = \operatorname{average}\{T_i\}$ ; and AM-Red-max, for which  $|F| = \max\{T_i\}$ .

Input data generation. Our experiments involved randomly generating task sets for a 16-core platform (so m=16) in three categories: light, medium, and heavy. Task utilizations ranged over (0.0,0.3) for light task sets, [0.3,0.7) for medium, and [0.7,1] for heavy. We generated masks independently of tasks. For hierarchical masks, we used a generation process that produced masks for which the number of nesting levels ranged from approximately m to approximately m. We discarded any non-feasible task sets that were generated.

**Experiments.** In our first experiment, we recorded, for each considered scheduler, the number of task migrations over an interval of length 10,000 time units as a function of relative system utilization, assuming periodic job releases. The *relative system utilization* of a task set is its total utilization divided by the number of cores. Task migrations are interesting to consider because semi-partitioned schedulers seek to limit them. Results can be found in Fig. 8. Compared to HPA-EDF, the curves shown here for the AM-Red variants are largely unaffected by the range of task utilizations (compare insets (a), (b), and (c) of Fig. 8). In contrast, HPA-EDF exhibits many more migrations for light task sets, which tend to have more tasks. Also, the curves for AM-Red-avg and AM-Red-max are significantly lower than those for HPA-EDF. Those for AM-Red-min start out lower, but in insets (b) and (c), eventually cross and become higher as relative system utilization nears one.

In our second experiment, we considered schedules similarly as above, but this time



**Figure 9** Exp. 2 (hierarchical masks): maximum tardiness, averaged over the generated task sets, as a function of relative system utilization.

computed maximum observed tardiness. Results can be found in Fig. 9. The curves shown here for the AM-Red variants are largely unaffected by the range of task utilizations (compare insets (a), (b), and (c) of Fig. 9). Also, each of these curves plateaus and remains steady beyond a certain relative system utilization value. In contrast, the curve for HPA-EDF is higher for larger task utilizations. Also, each HPA-EDF curve sharply increases as relative system utilization nears one. These curves indicate that tardiness tends to be predictably low only under AM-Red-min. However, from Fig. 8, this feature comes at the expense of more migrations in the case of medium and heavy task sets with high utilizations, which is not surprising.

We conclude this section by reminding the reader that HPA-EDF works only for hierarchical masks and has no schedulability test.

#### 8 Conclusion

We have presented AM-Red, the first (non-clairvoyant) optimal scheduler for implicit-deadline sporadic task sets assuming arbitrary processor affinity masks. We showed that AM-Red is SRT-optimal, with a tardiness bound of |F|, and that it is HRT-optimal if |F| divides the smallest task period. In addition, we presented analysis concerning task-migration frequency and time complexity. In the special case of hierarchical masks, we showed that AM-Red can be refined to find a valid frame in O(m+n) time, which is asymptotically optimal.

In other work that is omitted here due to space constraints, we have shown that the time complexity for frame construction can be reduced in other special cases. For example, for loop-free graphs, it can be reduced to O(m+n), which again is asymptotically optimal, using techniques similar to those discussed in Sec. 6. If masks are restricted in length, then it can also be reduced due to the internal structure of the affinity graph.

In future work, we intend to adapt AM-Red for heterogeneous multiprocessors, which are becoming more common in practice. Also, although the number of migrating tasks under AM-Red is generally optimal (i.e., task sets exist for which m-1 migrating tasks are fundamental, matching the bound in Theorem 18), we wish to find a way to reduce the number of migrating tasks to within a constant factor of that optimally required for each specific task set under consideration. Finally, while our focus in this paper has been semi-partitioned scheduling, global scheduling warrants consideration.

#### References

- QNX documentation. URL: http://www.qnx.com/developers/docs/qnxcar2/index.jsp?topic=%2Fcom.qnx.doc.multicore.user\_guide%2Ftopic%2Fhow\_to\_Thread\_affinity.html.
- 2 VxWorks 6.6 SMP. URL: http://archive.eettaiwan.com/www.eettaiwan.com/STATIC/PDF/200809/EETOL\_2008IIC\_WindRiver\_AN\_64.pdf?%20SOURCES=DOWNLOAD.
- 3 L. Abeni. SCHED\_DEADLINE: A real-time CPU scheduler for Linux. TuToR of Real-time Systems Symposium (RTSS 2017), 2017. URL: https://tutor2017.inria.fr/sched\_deadline/.
- 4 J. Anderson, V. Bud, and U. Devi. An EDF-based scheduling algorithm for multiprocessor soft real-time systems. In *Proceedings of Euromicro Conference on Real-Time Systems* (ECRTS 2005), pages 199–208. IEEE, 2005.
- 5 J. Anderson, J. Erickson, U. Devi, and B. Casses. Optimal semi-partitioned scheduling in soft real-time systems. *Journal of Signal Processing Systems*, 84(1):3–23, 2016.
- 6 S. Baruah and B. Brandenburg. Multiprocessor feasibility analysis of recurrent task systems with specified processor affinities. In *Proceedings of Real-Time Systems Symposium (RTSS 2013)*, pages 160–169. IEEE, 2013.
- 7 V. Bonifaci, B. Brandenburg, G. D'Angelo, and A. Marchetti-Spaccamela. Multiprocessor real-time scheduling with hierarchical processor affinities. In *Proceedings of Euromicro Conference on Real-Time Systems (ECRTS 2016)*, pages 237–247. IEEE, 2016.
- 8 V. Bonifaci, G. D'Angelo, and A. Marchetti-Spaccamela. Algorithms for hierarchical and semi-partitioned parallel scheduling. In *Proceedings of Parallel and Distributed Processing Symposium (IPDPS 2017)*, pages 738–747. IEEE, 2017.
- **9** D. Faggioli, F. Checconi, M. Trimarchi, and C. Scordino. An EDF scheduling class for the linux kernel. In *Proceedings of the Real-Time Linux Workshop (RTLWS 2009)*. IEEE, 2009.
- A. Foong, J. Fung, and D. Newell. An in-depth analysis of the impact of processor affinity on network performance. In *Proceedings of International Conference on Networks (ICN* 2004), volume 1, pages 244–250. IEEE, 2004.
- 11 L. Ford and D. Fulkerson. Maximal flow through a network. Canadian Journal of Mathematics, 8(3):399–404, 1956.
- P. Hall. On representatives of subsets. Journal of the London Mathematical Society, 1(1):26–30, 1935.
- 13 S. Kato, N. Yamasaki, and Y. Ishikawa. Semi-partitioned scheduling of sporadic task systems on multiprocessors. In *Proceedings of Euromicro Conference on Real-Time Systems* (ECRTS 2009), pages 249–258. IEEE, 2009.
- Y. Lee and A. Sidford. Path finding methods for linear programming: Solving linear programs in  $O(\sqrt{rank})$  iterations and faster algorithms for maximum flow. In *Proceedings of Foundations of Computer Science (FOCS 2014)*, pages 424–433. IEEE, 2014.
- J. Mistry, M. Naylor, and J. Woodcock. Adapting FreeRTOS for multicores: An experience report. Software: Practice and Experience, 44(9):1129–1154, 2014.
- 16 A. Mok. Fundamental Design Problems of Distributed Systems for Hard Real-Time Environments. PhD thesis, Massachusetts Institute of Technology, Cambridge, Mass., 1983.
- 17 G. Muneeswari and K. Shunmuganathan. A novel hard-soft processor affinity scheduling for multicore architecture using multi-agents. European Journal of Scientific Research, 55(3):419–429, 2011.
- 18 A. Ortiz, J. Ortega, A. Díaz, and A. Prieto. Affinity-based network interfaces for efficient communication on multicore architectures. *Journal of Computer Science and Technology*, 28(3):508–524, 2013.

#### YY:22 An Optimal Semi-Partitioned Scheduler Assuming Arbitrary Affinity Masks

- 19 K. Ramamritham and A. Stankovic. Scheduling algorithms and operating systems support for real-time systems. In *Proceedings of the IEEE*, volume 82, pages 55–67. IEEE, 1994.
- **20** A. Schrijver. *Combinatorial optimization: Polyhedra and efficiency*. Springer Science & Business Media, 2003.
- 21 M. Shekhar, A. Sarkar, H. Ramaprasad, and F. Mueller. Semi-partitioned hard-real-time scheduling under locked cache migration in multicore systems. In *Proceedings of Euromicro Conference on Real-Time Systems (ECRTS 2012)*, pages 331–340. IEEE, 2012.
- **22** J. Stefaniak and P. Wilson. System and method for assigning processes to specific CPU's to increase scalability and performance of operating systems, December 2 2003. US Patent 6,658,448.
- V. Williams. Multiplying matrices faster than Coppersmith-Winograd. In *Proceedings of Symposium on Theory of Computing (STOC 2012)*. ACM, 2012.
- K. Yang and J. Anderson. An optimal semi-partitioned scheduler for uniform heterogeneous multiprocessors. In *Proceedings of Real-Time Systems (ECRTS 2015)*, pages 199–210. IEEE, 2015.
- P. Zijlstra. An update on Real-Time scheduling on Linux. Keynote talk at Euromicro Conference on Real-Time Systems (ECRTS 2017), 2017. URL: https://www.ecrts.org/fileadmin/files\_ecrts17/ecrts17-peterz.pdf.