

Algorithms for Hierarchical and Semi-Partitioned Parallel Scheduling

Vincenzo Bonifaci^{a,*}, Gianlorenzo D'Angelo^b, Alberto Marchetti-Spaccamela^c

^a*IASI-CNR, Rome, Italy*

^b*Gran Sasso Science Institute, L'Aquila, Italy*

^c*Sapienza Università di Roma, Rome, Italy*

Abstract

We propose a model for scheduling jobs in a parallel machine setting that takes into account the cost of migrations by assuming that the processing time of a job may depend on the specific set of machines among which the job is migrated. For the makespan minimization objective, the model generalizes classical scheduling problems such as unrelated parallel machine scheduling, as well as novel ones such as semi-partitioned and clustered scheduling. In the case of a hierarchical family of machines, we derive a compact integer linear programming formulation of the problem and leverage its fractional relaxation to obtain a polynomial-time 2-approximation algorithm. Extensions that incorporate memory capacity constraints are also discussed.

Keywords: processor affinities, makespan minimization, unrelated machines, laminar family, wrap-around rule, clustered scheduling

1. Introduction

Multicore architectures have become the standard computing platform in many domains: multicore processors speed up application performance by dividing the workload among multiple processing cores instead of using one “super fast” single processor. A hierarchical organization of chips or clusters of symmetric multiprocessing (SMP) nodes with multicore chip-multiprocessors (CMP), also known as SMP-CMP clusters, is common today. For example, consider the architecture of Intel’s Dual-Core Xeon (see Figure 1). In this architecture there are three levels of communication: the communication between two processors on the same chip (intra-CMP communication); the communication across chips but within a node (inter-CMP communication), and the communication

*Corresponding author

Email addresses: vincenzo.bonifaci@iasi.cnr.it (Vincenzo Bonifaci), gianlorenzo.dangelo@gssi.it (Gianlorenzo D'Angelo), alberto@diag.uniroma1.it (Alberto Marchetti-Spaccamela)

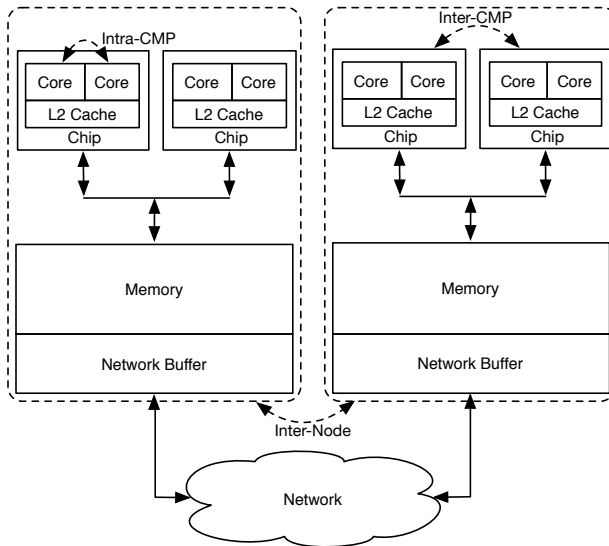


Figure 1: Example of a multicore cluster.

between two processors on different nodes (inter-node communication). Intra-CMP communication has higher performance than inter-CMP, which in turn has higher performance than inter-node communication: communications between cores within a chip can usually achieve lower latency and higher bandwidth than communications between cores in different chips.

The objective of how to efficiently exploit the available hardware parallelism for scheduling jobs is crucial. Experimental work —see for example [26] and references therein— has shown that a dynamic scheduler that tries to balance the processes among the available resources to ensure fair distribution of CPU time and minimize idle cores is not sufficient. The fundamental flaw with this approach is that a core is not an independent processor, but is part of a larger on-chip system that shares resources (such as caches and buses) with other cores. For example, in the multicore system of the dual-core Xeon, cores on the same chip share the same L2 cache and memory controller, and all the cores access the main memory through a shared bus.

Since the communication cost is not uniform, the cost of preempting a job and resuming its execution should take into account the involved cores: the cost of resuming execution of a job on the same core is lower than the cost of resuming on a different core; moreover, the cost of migration is not uniform and depends on the communication cost between the two involved cores. For all these reasons, scheduling policies are needed that not only limit the number of migrations, but that are also aware of the costs involved in migrations.

Additionally, we note that there is a trend in the design of multicore architectures towards heterogenous architectures, providing more flexibility to meet

specific performance/energy consumption goals. In fact, heterogeneous multi-core architectures have been shown to require significantly less energy without a significant degradation of performance. This results in higher overall efficiency with respect to conventional homogeneous architectures, but implies that the processing time of a job cannot be regarded as a constant.

In this paper, we propose a theoretical model for scheduling jobs in a multi-core architecture that can capture the cost of migrations by assuming that *the processing time of a job depends on the specific set of machines on which the job is scheduled*. Namely, we are given a family of admissible sets of machines \mathcal{A} , and, for each job j and for each set $\alpha \in \mathcal{A}$, a value $P_j(\alpha)$ denoting the processing time required by j if its execution is limited to the machines in α . We assume that jobs that are assigned to a set α can be executed on any machine in the set; they can be preempted and possibly migrated to another machine in α , but simultaneous processing of the same job by two machines is not allowed (see Section 2 for details).

This setting opens up a whole new class of scheduling models with their own particular challenges and subsumes well-known problems. For example, if there are m machines and the admissible family \mathcal{A} consists of singletons (i.e., $\mathcal{A} = \{\{1\}, \{2\}, \dots, \{m\}\}$), then we obtain the unrelated machines scheduling problem [15]; if \mathcal{A} consists of one set containing all machines (i.e., $\mathcal{A} = \{\{1, 2, \dots, m\}\}$) then we have the (preemptive) parallel machine scheduling problem [20].

While the model presented here does not account exactly for the number of migrations incurred, this number can be bounded (see, for example, Proposition 3.4), allowing migration costs to be accounted for in the processing times, if desired. Differently from other approaches, this allows for a flexible input representation and easily accommodates heterogeneous machines. In fact, we also show how our basic model can be extended to incorporate memory capacity constraints.

1.1. Related Work

Much of the prior work on multiprocessor scheduling theory has focused on either the *partitioned* or the *global* approach. Under partitioning, each job is statically assigned to a machine; if the cost of processing a job depends on the specific machine on which the job is executed, we have the unrelated machines scheduling problem [15]. Under global scheduling, on the other hand, task migration is allowed with no restrictions and with no additional costs [14, 20].

It is well-known that partitioning incurs lower runtime overheads (as there are no migrations), but produces schedules that may be unnecessarily constrained; global scheduling, vice versa, entails higher runtime costs that should be properly taken into account (see for example [25]). We now review other approaches that have been proposed and experimentally tested to overcome the above tradeoff between better scheduling policies and higher costs.

Semi-partitioned scheduling was proposed as a compromise between pure partitioned and global scheduling [3]. Semi-partitioning relaxes partitioned scheduling by allowing a small number of jobs to migrate, thereby improving

schedulability. Such tasks are called migratory, in contrast to statically assigned tasks. The common goal in this line of work is to circumvent the algorithmic limitations and resulting capacity loss of partitioning, while avoiding the overhead of global scheduling by limiting migrations.

Clustered scheduling is another proposal that aims to alleviate limitations of partitioned and global algorithms; it exploits the grouping of cores into clusters of symmetric multiprocessing nodes with multicore chip multiprocessors: tasks are statically assigned to clusters (like in partitioning), but are globally scheduled within each cluster [2, 22].

Semi-partitioned and clustered scheduling are not the only two proposals; we briefly mention other proposals. *Federated scheduling* was introduced in [16] to deal with parallel real-time tasks, where each task is a DAG whose nodes represent jobs and edges represent precedence constraints among jobs. Forcing the execution of a single task on a single processor restricts all jobs of a task to execute on the same processor, and forbids to deal with tasks with a (parallelizable) computational demand exceeding the capacity of a single processor. The federated scheduling approach [1, 16] is advocated as a reasonable extension of partitioned scheduling to parallel real-time task sets: there are tasks that are permitted to execute upon more than one processor and are granted exclusive access to the processors upon which they execute, while the remaining tasks are partitioned amongst a pool of shared processors.

We observe that contemporary commodity operating systems (such as Linux and Windows) implement more complex migration strategies by defining *processor affinity masks*, which specify on a per-process basis on which processors a job may be scheduled. Namely, processor affinities allow binding a process to an arbitrary subset of processors in the system and a process can only be scheduled on the processors that it is bound to. It is known that processor affinity is useful for increasing the performance of a parallel system in several contexts, such as application performance, fault-tolerance, security and real-time systems [5, 19, 23].

Processor affinities yield a general framework that can be used to realize global, partitioned, or clustered scheduling. For example, in partitioned scheduling, each task’s processor affinity includes exactly one processor, while in global scheduling, each task’s processor affinity is set to all processors. The new feature is that arbitrary processor affinities can be assigned on a job-by-job basis, which permits the specification of migration strategies that are more flexible than those usually studied in the literature.

To the best of our knowledge, the model we propose has not been considered theoretically. In the rest of this section we highlight the differences with previous results on somewhat similar, but distinct, models.

We already observed that the problem of scheduling unrelated machines is a special case of our model; more recently, in [8, 21] a non-preemptive scheduling problem is considered that is a special case of scheduling unrelated machines (and, thus, of our model). Namely, the problem of non-preemptively scheduling to minimize makespan when each job j can be scheduled only on a set of machines $M(j)$ (i.e., the processing time of job j on machine i is p_j if $i \in M(j)$)

or ∞ otherwise); such machine sets have a laminar structure. Glass and Kellerer [8] give a $(2 - 1/m)$ -approximation algorithm; Muratore et al. [21] improve this to a polynomial-time approximation scheme.

Bougeret et al. [6] (see also references therein) consider the non-preemptive scheduling of jobs on a clustered architecture. The authors assume that each cluster is formed by m processors and that the execution of job j requires q_j processors belonging to the same cluster for p_j time units; they develop a $7/3$ approximation algorithm for minimizing the makespan. In our model we allow preemption and we consider the more challenging case when for each job there are *many* sets of machines that could execute it, with different processing times; indeed, one of the main difficulties lies in selecting the best processor affinity mask for each job.

Hwang et al. [10] study a model of parallel non-preemptive scheduling on identical machines, where interprocessor communication times are explicitly given as part of the input. While potentially more accurate, we observe that in order to be applied to a preemptive setting, such a model would require to break down each job into a possibly exponential number of unit-size jobs; additionally, the model of Hwang et al. would require a significant extension in the heterogeneous case.

Our work focuses on the cases of hierarchical affinity masks. In the case of general, non-hierarchical affinity masks with no memory capacity constraints, one can obtain an 4-approximation algorithm by invoking known results, as follows. Given an instance of the general problem, construct an instance of the unrelated machines problem by setting the processing time of a job j on machine i to the minimum processing time of j on an admissible set that contains i . The optimal *preemptive* makespan of such an instance is a lower bound on the optimum of the original instance with affinity masks. Moreover, Lin and Vitter [17] (see also [7]) show how to construct in polynomial time a *non-preemptive* schedule whose makespan is at most a factor 4 times the optimal preemptive makespan. Such a schedule is then a 4-approximate solution for the problem of scheduling with affinity masks.

Finally, we remark that in this work we focus on the load balancing and runtime scheduling aspects rather than memory accesses and cache complexity. The reader is referred to [4] and references therein for models that focus on hierarchical cache performance.

1.2. Contributions and Organization of the Paper

We focus on preemptively scheduling jobs to minimize the makespan assuming a hierarchical architecture, and we first show how this problem generalizes several classical and new problems (Section 2).

In Section 3 we consider, as a warm-up, the case of semi-partitioned scheduling, and we identify necessary and sufficient conditions for schedulability. Namely, we provide an ILP formulation of the assignment problem that, for each job, will specify whether the job is assigned to a specific machine or executed globally. We also provide an efficient scheduler that, given a feasible solution to the

ILP, constructs a schedule with the same makespan, thus setting the times for executing and possibly migrating jobs.

In Section 4 we consider the more general case of hierarchical scheduling. We provide an ILP formulation of the assignment problem that, for each job, will specify the *affinity mask* that will be used for scheduling the job, and a scheduler that, given a feasible solution to the ILP, constructs a schedule with the same makespan, again setting the times for executing and possibly migrating jobs. We remark that our proposed scheduler is combinatorial and that the schedule cannot be trivially constructed by using standard network flow formulations for scheduling on identical machines.

In Section 5 we prove an upper bound on the approximation ratio of the problem for the hierarchical setting. We show how a fractional relaxation of the ILP can be leveraged to obtain a polynomial-time 2-approximation algorithm by building on existing algorithms for the unrelated machines scheduling problem; the key lemma of the proof shows how to redistribute the variables' values in a feasible fractional solution across the levels of the machine hierarchy.

Finally, in Section 6 we consider extensions of the model to handle additional memory constraints: each job is characterized by a memory requirement in addition to its processing times, and there is a constraint on the available memory at each machine, or at each cluster of the hierarchy.

2. Notation and problem formulation

We are given a set of n jobs $J := \{1, \dots, n\}$ and a set of m machines $M := \{1, \dots, m\}$. Each job needs to be assigned to a *set* of machines on which the job is allowed to be schedule, and the job can be preempted and migrated among any such machines. However, its processing time depends on the set of machines on which it is assigned. In detail, we are given a family of admissible sets $\mathcal{A} \subseteq 2^M$, and for each job $j \in J$, a processing time function $P_j : \mathcal{A} \rightarrow \mathbb{Z}_+$ with the constraint that the function must be *monotone* on \mathcal{A} , i.e., if $\alpha, \beta \in \mathcal{A}$ and $\alpha \subseteq \beta$, then $P_j(\alpha) \leq P_j(\beta)$, modeling the fact that processing overheads (caused, e.g., by migration) increase if the job is executed using a larger set of machines. We often use the shorthand $p_{\alpha j} := P_j(\alpha)$. Moreover, to avoid cumbersome notation, when α is a singleton, such as $\alpha = \{i\}$, we write p_{ij} instead of $p_{\{i\}j}$.

We therefore stipulate that, when a job $j \in J$ is run on a set of machines $\alpha \in \mathcal{A}$, the total processing time it receives must be $P_j(\alpha)$. In general, if a job is run on machine set M' (which may or may not be in \mathcal{A}), its processing time must be $p_{\alpha j}$, where α is an inclusion-wise minimal set in \mathcal{A} that contains M' and that minimizes the processing time (if there is no such α , then j cannot be run on M').

Given J and \mathcal{A} , an *assignment* of jobs in J to sets in \mathcal{A} is a function that assigns each job in J to a set in \mathcal{A} . If a job j is assigned to a set α , then its *processing time* is $P_j(\alpha)$. The set α to which a job j is assigned is also called the *affinity mask* of j . Given an assignment of jobs in J to sets in \mathcal{A} , a schedule

is *valid* with respect to the assignment if each job is scheduled on time slots of machines in its affinity mask, no job is processed in parallel on more than one machine in the same time interval (though it may be preempted or migrated), each job receives the required amount of processing time (i.e. $p_{\alpha j}$, if job j is assigned to set α), and no machine processes more than one job in a time slot. We assume that schedules start at time 0 and allow preemptions and migrations to occur only at integer time points. If, in a given schedule, a job j completes at time C_j , then $T := \max_{j \in J} C_j$ is called the *makespan* of the schedule.

In this paper, we consider the problem of finding an assignment of jobs in J to sets in \mathcal{A} and a corresponding valid schedule that minimizes the makespan. We divide the problem into two subproblems: given J and \mathcal{A} , find an assignment of jobs in J to sets in \mathcal{A} that admits a valid schedule in the interval $[0, T]$ and minimizes T ; and given an assignment of jobs in J to sets in \mathcal{A} that admits some valid schedule in the interval $[0, T]$, construct a valid schedule in the same interval.

In this paper, we restrict the discussion to *laminar* (or *hierarchical*) instances of the problem, where, for each $\alpha, \alpha' \in \mathcal{A}$, either $\alpha \subseteq \alpha'$ or $\alpha' \subseteq \alpha$ or $\alpha \cap \alpha' = \emptyset$. Without loss of generality, we assume that all sets in the family \mathcal{A} are distinct. In a laminar instance, the *level* of a set β is the number of sets $\alpha \in \mathcal{A}$ such that $\beta \subseteq \alpha$ and the level of the instance is the maximum level among the sets in \mathcal{A} . We call the problem with laminar instances the *hierarchical scheduling problem*. The hierarchical scheduling problem generalizes some well-known and new scheduling problems.

- *Identical parallel machines* scheduling with preemption ($P|pmtn|C_{\max}$) [20]: take $\mathcal{A} = \{M\}$. Then each job j can be migrated freely among the machines in M , as long as it receives the processing time p_{Mj} .
- *Unrelated parallel machines* scheduling ($R||C_{\max}$) [15]: take \mathcal{A} to be a family of m singletons, i.e., $\mathcal{A} = \{\{1\}, \{2\}, \dots, \{m\}\}$. Then each job must be assigned to a single machine (no migration) and its processing time is a function of the machine.
- *Semi-partitioned* scheduling [3]: take $\mathcal{A} = \{M, \{1\}, \{2\}, \dots, \{m\}\}$. Then each job can either be run *globally* (i.e., freely migrated) on M with processing time p_{Mj} , or assigned *locally* to a specific machine $i \in M$, with processing time $p_{ij} \leq p_{Mj}$.
- *Clustered* scheduling [2]: let $m = kq$. Take $\mathcal{A} = \{M, \{1\}, \dots, \{m\}, \{1, \dots, q\}, \{q+1, \dots, 2q\}, \dots, \{(k-1)q, \dots, kq\}\}$. Then each job can be run globally, or locally to a single machine, or locally to a *cluster* of q machines.

Semi-partitioned scheduling generalizes scheduling on unrelated parallel machines (by taking sufficiently large values of p_{Mj}); hence, the following proposition is implied by existing results for the $R||C_{\max}$ problem [15].

Proposition 2.1. *Hierarchical and semi-partitioned scheduling are NP-hard to approximate within any constant factor less than $3/2$.*

The following example shows that, not surprisingly, hierarchical scheduling instances may admit shorter schedules than the corresponding unrelated parallel machine instances.

Example 2.1. (See also Figure 2.) Consider a semi-partitioned instance with three jobs and two machines: job 1 has $p_{M1} = \infty$, $p_{11} = 1$, $p_{21} = \infty$; job 2 has $p_{M2} = \infty$, $p_{12} = \infty$, $p_{22} = 1$; job 3 has $p_{M3} = p_{13} = p_{23} = 2$ (∞ represents a sufficiently large constant). It is easy to see that the semi-partitioned instance has a schedule with makespan 2, while the corresponding unrelated machine instance has an optimal makespan of 3.

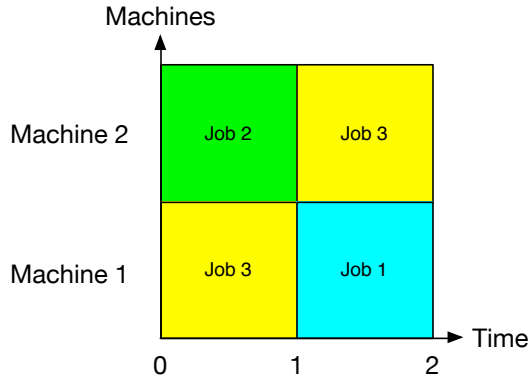


Figure 2: Optimal schedule for the instance of Example 2.1.

In the sequel, to more easily illustrate the ideas behind our approach, we first discuss the two-level case in Section 3. Section 4 is devoted to the more general hierarchical setting. We describe the details of the rounding procedure in Section 5, which leads to our main result (Theorem 5.2). Finally, in Section 6, we discuss how the model can be extended to incorporate memory capacity constraints.

3. Semi-partitioned scheduling

Let $j = 1, \dots, n$ be a job and $i = 1, \dots, m$ be a machine; in this section, the special “machine” index 0 will represent the set M , i.e., global processing. We use the following integer linear program (ILP) to determine the minimum makespan. Binary variable x_{ij} encodes the assignment of job j to machine i (or

to the set M , if $i = 0$).

$$\min T \tag{IP-1}$$

$$\sum_{i=0}^m x_{ij} = 1 \quad \text{for } j \in J \tag{1a}$$

$$\sum_{j=1}^n \sum_{i=0}^m p_{ij} x_{ij} \leq mT \tag{1b}$$

$$\sum_{j=1}^n p_{ij} x_{ij} \leq T \quad \text{for } i \in M \tag{1c}$$

$$p_{ij} x_{ij} \leq T \quad \text{for } j \in J \text{ and } i \in \{0\} \cup M \tag{1d}$$

$$x_{ij} \in \{0, 1\} \quad \text{for } j \in J \text{ and } i \in \{0\} \cup M. \tag{1e}$$

It should be clear from the description of the model that the above constraints are necessary for the existence of a valid schedule with makespan T ; the fact that they are sufficient is the nontrivial claim that we prove in this section. We show algorithmically that a feasible solution to (IP-1) ensures that a valid schedule with makespan T exists.

Example 3.1. Consider the instance of Example 2.1. For any finite value of T , the ILP constraints imply $x_{11} = 1$, $x_{22} = 1$; thus, for job 3, we obtain the processing time constraints $2x_{03} \leq T$, $2x_{13} \leq T - 1$, $2x_{23} \leq T - 1$, $2x_{13} + 2x_{23} + 2x_{03} \leq 2T - 2$. The optimal integral solution has $T = 2$ and assigns job 1 to machine 1, job 2 to machine 2, and job 3 globally to both machines. The following schedule has a makespan value of 2: job 1 is scheduled on machine 1 during $[1, 2)$; job 2 is scheduled on machine 2 during $[0, 1)$; finally, job 3 is scheduled on machine 1 during $[0, 1)$, then migrated to machine 2 where it is scheduled during $[1, 2)$.

The pseudo-code of our scheduler is reported in Algorithm 1. The scheduler takes as input a feasible solution to (IP-1) and assigns the jobs to time slots of the machines according to the affinity masks defined in the solution.

Algorithm 1 first schedules the *global jobs* (i.e. jobs j such that $x_{0j} = 1$) so that no job is scheduled simultaneously on two machines. Namely, it computes the total volume V of global jobs and initializes a variable t to 0 (lines 1–2). Then, it iterates on all the machines and assigns to each of them a suitable amount δ of global volume. While $V > 0$, the algorithm looks for an empty machine $i > 0$ (line 4) and schedules δ global volume in the interval $[t, t + \delta \pmod{T}]$ on i (line 6). Then it increases the value of t by $\delta \pmod{T}$ and decreases the volume V of global jobs still to be scheduled by δ (lines 7–8).

The value of δ in each iteration is computed as follows. The total volume of *local jobs* assigned to machine i is $\sum_{j=1}^n p_{ij} x_{ij}$, so we can schedule at most $T - \sum_{j=1}^n p_{ij} x_{ij}$ volume of *global jobs* on i in the interval $[0, T]$. Therefore, if the volume V of global jobs that still needs to be scheduled is smaller than

ALGORITHM 1: Job scheduling (for a given assignment \mathbf{x})

```
1  $t \leftarrow 0$ ;  
2  $V \leftarrow \sum_{j=1}^n p_{0j}x_{0j}$ ;  
3 while  $V > 0$  do  
4    $i \leftarrow$  an empty machine (in  $1, \dots, m$ );  
5    $\delta \leftarrow \min(V, T - \sum_{j=1}^n p_{ij}x_{ij})$ ;  
6   Assign  $\delta$  units of global work to  $i$ , in the interval  $[t, t + \delta \pmod{T}]$ ;  
7   /* (The remaining  $T - \delta$  units on  $i$  will be used by local jobs) */  
8    $t \leftarrow t + \delta \pmod{T}$ ;  
9    $V \leftarrow V - \delta$ ;  
10 end  
11 foreach machine  $i \in M$  and job  $j \in J$  such that  $x_{ij} = 1$  do  
12 | Schedule  $j$  on machine  $i$  in the free time of interval  $[0, T]$ ;  
13 end
```

$T - \sum_{j=1}^n p_{ij}x_{ij}$, then $\delta = V$, otherwise, we exploit all the possible empty space on i and then $\delta = T - \sum_{j=1}^n p_{ij}x_{ij}$ (line 5).

Having scheduled the global jobs, the algorithm then schedules the *local* jobs (i.e. jobs j such that $x_{ij} = 1$, for $i > 0$) in the free time of each machine (line 11).

In order to show that Algorithm 1 produces a valid schedule, we need to prove the two next lemmas.

Lemma 3.1. *In the schedule produced by Algorithm 1, all jobs receive the required amount of processing time.*

Proof. We first show the statement for global jobs, in particular we show that $V = 0$ at the end of the while loop. By contradiction, assume that at the end of some iteration of the while loop $V > 0$ and there is no empty machine left to be selected at line 4 of the next iteration. This implies that for each machine i , the amount δ of global volume scheduled on i is $T - \sum_{j=1}^n p_{ij}x_{ij}$ (see lines 5 and 8). Therefore, the overall amount of scheduled global jobs is $\sum_{i=1}^m (T - \sum_{j=1}^n p_{ij}x_{ij})$ and the amount of global jobs still to be scheduled is

$$\sum_{j=1}^n p_{0j}x_{0j} - \sum_{i=1}^m (T - \sum_{j=1}^n p_{ij}x_{ij}) = \sum_{i=0}^m \sum_{j=1}^n p_{ij}x_{ij} - mT.$$

Since $V > 0$ at the end of the considered iteration, then the above quantity is strictly positive, and then $\sum_{i=0}^m \sum_{j=1}^n p_{ij}x_{ij} > mT$, a contradiction to constraint (1b).

Note also that, for each machine i , the global jobs leave a free time of at least $T - \sum_{j=1}^n p_{ij}x_{ij}$ in the interval $[0, T]$ (see line 5). Therefore, the local jobs that are assigned to machine i receive at least an overall amount $\sum_{j=1}^n p_{ij}x_{ij}$ of processing time at line 11. \square

Lemma 3.2. *In the schedule produced by Algorithm 1, no job is scheduled in parallel with itself.*

Proof. Clearly, no local job will be scheduled in parallel with itself, since each such job is scheduled on a unique machine (and $p_{ij}x_{ij} \leq T$). Assume by contradiction that a global job j is scheduled on two different machines i and i' during the same time interval $[t_1, t_2]$, $t_2 > t_1$, and assume w.l.o.g. that the iteration related to machine i occurred before that of i' . This implies that in all the iterations from that of i to that of i' , only job j is scheduled and that, since $t_2 > t_1$, $p_{0j} > T$. Since $x_{0j} = 1$ this is a contradiction to constraint (1d). \square

Theorem 3.3. *Given a feasible solution (\mathbf{x}, T) to (IP-1), Algorithm 1 produces a valid schedule in the interval $[0, T]$.*

Proof. The statement follows from lemmas 3.1 and 3.2 and by observing that all the jobs are scheduled in the interval $[0, T]$ (see lines 6–7 and 11 of Algorithm 1). \square

The following proposition bounds the number of migrations and preemptions that occur in the worst case; this value can be used for a priori bounding the worst-case processing time of a job that is migrated among multiple machines.

Proposition 3.4. *The number of job migrations in the schedule produced by Algorithm 1 is at most $m - 1$. The number of job preemptions and migrations is at most $2m - 2$.*

4. Hierarchical scheduling

The following ILP expresses necessary conditions on the minimum makespan T and an optimal assignment \mathbf{x} in the hierarchical scheduling problem.

$$\min T \tag{IP-2}$$

$$\sum_{\alpha \in \mathcal{A}} x_{\alpha j} = 1 \quad \text{for } j \in J \tag{2a}$$

$$\sum_{j=1}^n \sum_{\beta \subseteq \alpha} p_{\beta j} x_{\beta j} \leq |\alpha|T \quad \text{for } \alpha \in \mathcal{A} \tag{2b}$$

$$p_{\alpha j} x_{\alpha j} \leq T \quad \text{for } \alpha \in \mathcal{A}, j \in J \tag{2c}$$

$$x_{\alpha j} \in \{0, 1\} \quad \text{for } \alpha \in \mathcal{A}, j \in J. \tag{2d}$$

Similarly to the previous section, we give an algorithm that takes as input a feasible solution (\mathbf{x}, T) to (IP-2) and constructs a valid schedule with makespan T . The algorithm works in two phases. The first phase (Algorithm 2) proceeds bottom-up (i.e., from the smallest sets up to the largest set) and, for each set of machines $\alpha \in \mathcal{A}$ and for each machine $i \in \alpha$, it determines the load $\text{LOAD}[i, \alpha]$ of machine i due to jobs assigned to set α by the ILP (i.e., jobs j

s.t. $x_{\alpha j} = 1$). The load is assigned in such a way that for each affinity mask the overall load is equal to the sum of the required processing time, that is $\sum_{i \in \alpha} \text{LOAD}[i, \alpha] = \sum_{j=1}^n p_{\alpha j} x_{\alpha j}$.¹ The second phase (Algorithm 3) proceeds top-down (i.e., from the largest set down to the smallest ones) and, for each set $\alpha \in \mathcal{A}$, determines the schedule of each job j such that $x_{\alpha j} = 1$ on each machine $i \in \alpha$, that is, the time slots of each machine $i \in \alpha$ that are assigned to job j .

The crucial observation is that the first phase computes LOAD in such a way that the second phase is able to identify only one machine, for each affinity mask, that must be checked in order to avoid to schedule more than one job in the same time interval of the same machine. In fact, the first phase ensures that for each affinity mask $\beta \in \mathcal{A}$ there exists at most one machine $i \in \beta$ that is loaded with some jobs assigned to a superset of β , that is $\text{LOAD}[i, \beta] > 0$ and $\text{LOAD}[i, \alpha] > 0$, for some $\alpha \in \mathcal{A}$ such that $\beta \subset \alpha$. Since the second phase proceeds top-down, when affinity mask β is analyzed, the schedule of jobs assigned to α in such a machine i is already determined.

Assume that the jobs assigned to α are scheduled in interval $[t, t_{i\alpha}]$, where $t_{i\alpha} = t + \text{LOAD}[i, \alpha] \pmod{T}$; the algorithm first schedules the jobs assigned to β in machine i , in the interval $[t_{i\alpha}, t_{i\beta}]$, where $t_{i\beta} = t_{i\alpha} + \text{LOAD}[i, \beta] \pmod{T}$. Then, it schedules the remaining load by filling up machines $\ell \in \beta \setminus \{i\}$ starting from time $t_{i\beta}$. Indeed, if we assume that the machines in $\beta \setminus \{i\}$ are sorted in an arbitrary way, $\beta \setminus \{i\} = (\ell_1, \ell_2, \dots, \ell_{|\beta|-1})$, then jobs assigned to β are scheduled in interval $[t_{\ell_{k-1}\beta}, t_{\ell_k\beta}]$ of machine ℓ_k , where $t_{\ell_0\beta} = t_{i\beta}$ and $t_{\ell_k\beta} = t_{\ell_{k-1}\beta} + \text{LOAD}[\ell_k, \beta] \pmod{T}$. This guarantees that no jobs is scheduled in parallel with itself and that no machine has more than one job scheduled in the same time interval.

The pseudo-code of the first phase is given in Algorithm 2. First, the algorithm initializes variable LOAD (line 1). Then, for each $\alpha \in \mathcal{A}$ and $i \in \alpha$, it initializes variable $\text{TOT-LOAD}[i, \alpha]$, which stores the cumulative load of machine i due to all sets $\beta \subseteq \alpha$ (line 2). Variable $\text{MARK}[\alpha]$ is used to determine whether set $\alpha \in \mathcal{A}$ has been visited or not by the algorithm and it is initialized at line 3. The while loop at lines 4–14 visits all the sets in \mathcal{A} in a bottom-up order: at each iteration it selects a set α such that all its subsets have been already visited, i.e. such that $\text{MARK}[\alpha] = \text{false}$ and, for each $\beta \subset \alpha$, $\text{MARK}[\beta] = \text{true}$ (line 5). At each iteration, variable V stores the volume of jobs assigned to α (i.e. jobs j such that $x_{\alpha j} = 1$) that still needs to be scheduled. Variable V is initialized to the total volume of jobs assigned to α at line 6. The loop at lines 7–13 iterates for each machine $i \in \alpha$ in ascending order and, in order to compute $\text{LOAD}[i, \alpha]$, first selects the maximal subset β of α that contains machine i (if it exists, see line 8–

¹An alternative approach could be to modify (IP-2) by adding additional fractional variables of the form $y_{\alpha ij}$ and constraints of the form $\sum_i y_{\alpha ij} = x_{\alpha j}$; $y_{\alpha ij}$ represents the fractional share of job j on machine i if job j is scheduled using affinity mask α . However, this may not suffice to guarantee that a valid schedule exists, since a job can only be scheduled on one machine at a time. Conversely, our approach guarantees that a valid schedule exists (Theorem 4.6); moreover, the method is combinatorial and avoids the complication of a larger number of variables.

ALGORITHM 2: First phase (bottom-up volume allocation)

```
1 LOAD[ $i, \alpha$ ] = 0, for each  $\alpha \in \mathcal{A}$  and  $i \in \alpha$ ;
2 TOT-LOAD[ $i, \alpha$ ] = 0,  $\alpha \in \mathcal{A}$  and  $i \in \alpha$ ;
3 MARK[ $\alpha$ ] = false for each  $\alpha \in \mathcal{A}$ ;
4 while  $\exists \alpha \in \mathcal{A}$  such that  $\neg$ MARK[ $\alpha$ ] do
5   Let  $\alpha$  such that  $\neg$ MARK[ $\alpha$ ] and (MARK[ $\beta$ ], for each  $\beta \subset \alpha$ );
6    $V \leftarrow \sum_{j=1}^n p_{\alpha j} x_{\alpha j}$ ;
7   foreach  $i \in \alpha$  in ascending order do
8     Let  $\beta$  be the maximal set  $\beta \subset \alpha$  such that  $i \in \beta$ ;
9     (if no such  $\beta$  exists, set  $\beta = \emptyset$  and TOT-LOAD[ $i, \emptyset$ ] = 0);
10    LOAD[ $i, \alpha$ ]  $\leftarrow \min\{V, T - \text{TOT-LOAD}[i, \beta]\}$ ;
11    TOT-LOAD[ $i, \alpha$ ]  $\leftarrow \text{TOT-LOAD}[i, \beta] + \text{LOAD}[i, \alpha]$ ;
12     $V \leftarrow V - \text{LOAD}[i, \alpha]$ ;
13  end
14  MARK[ $\alpha$ ]  $\leftarrow$  true;
15 end
```

9). The value of $\text{LOAD}[i, \alpha]$ is computed at line 10 as follows. The total volume of jobs already assigned to machine i is equal to the cumulative load of machine i due to all sets $\beta \subset \alpha$, that is $\text{TOT-LOAD}[i, \beta]$. Then, we can schedule at most $T - \text{TOT-LOAD}[i, \beta]$ volume of jobs assigned to α on i . Therefore, if the volume V of global jobs that still needs to be scheduled is smaller than $T - \text{TOT-LOAD}[i, \beta]$, then we assign the entire volume to i and set $\text{LOAD}[i, \alpha] = V$, otherwise, we exploit all the possible empty space and set $\text{LOAD}[i, \alpha] = T - \text{TOT-LOAD}[i, \beta]$. Next, the algorithm computes the value of $\text{TOT-LOAD}[i, \alpha]$ by adding $\text{LOAD}[i, \alpha]$ to $\text{TOT-LOAD}[i, \beta]$ (line 11). Note that, $\text{TOT-LOAD}[i, \alpha] = \sum_{\beta \subset \alpha: i \in \beta} \text{LOAD}[i, \beta]$ and, eventually, $\sum_{i \in \alpha} \text{TOT-LOAD}[i, \alpha] = \sum_{\beta \subset \alpha} \sum_{i \in \beta} \text{LOAD}[i, \beta]$. Finally, variables V and $\text{MARK}[\alpha]$ are updated at lines 12 and 14. The next lemma shows that the cumulative load on each machine i is at most T and that the volume of the jobs assigned to set α is assigned entirely to variables $\text{LOAD}[i, \alpha]$, for all $i \in \alpha$.

Lemma 4.1. i) For every $\alpha \in \mathcal{A}$ and $i \in \alpha$, $\text{TOT-LOAD}[i, \alpha] \leq T$ at the end of Algorithm 2.

ii) Whenever line 14 of Algorithm 2 is executed, $V = 0$.

Proof. i) From lines 10 and 11 of the algorithm we get $\text{TOT-LOAD}[i, \alpha] \leq \text{TOT-LOAD}[i, \beta] + T - \text{TOT-LOAD}[i, \beta] = T$. ii) Let α be the first set for which the statement does not hold. That is, $V > 0$ at the end of line 13 of Algorithm 2 of the iteration related to set α , while $V = 0$ at the end of the iteration related to each $\beta \subset \alpha$.

For each $\beta \subset \alpha$, we have that

$$\sum_{i \in \beta} \text{LOAD}[i, \beta] = \sum_{j=1}^n p_{\beta j} x_{\beta j},$$

since $V = \sum_{j=1}^n p_{\beta j} x_{\beta j} - \sum_{i \in \beta} \text{LOAD}[i, \beta] = 0$. By definition of TOT-LOAD, for each $\beta \subset \alpha$, $\sum_{i \in \beta} \text{TOT-LOAD}[i, \beta] = \sum_{\gamma \subseteq \beta} \sum_{i \in \gamma} \text{LOAD}[i, \gamma]$ and then,

$$\sum_{i \in \beta} \text{TOT-LOAD}[i, \beta] = \sum_{\gamma \subseteq \beta} \sum_{j=1}^n p_{\gamma j} x_{\gamma j}. \quad (3)$$

For each $i \in \alpha$, let β_i be the maximal set $\beta_i \subset \alpha$ such that $i \in \beta_i$ (see line 8). Since for α the statement does not hold, then, for each $i \in \alpha$, $\text{LOAD}[i, \alpha] = \min\{V, T - \text{TOT-LOAD}[i, \beta_i]\} = T - \text{TOT-LOAD}[i, \beta_i]$ (see line 10). It follows that at line 14 of the iteration related to set α :

$$\begin{aligned} V &= \sum_{j=1}^n p_{\alpha j} x_{\alpha j} - \sum_{i \in \alpha} \text{LOAD}[i, \alpha] \\ &= \sum_{j=1}^n p_{\alpha j} x_{\alpha j} - \sum_{i \in \alpha} T + \sum_{i \in \alpha} \text{TOT-LOAD}[i, \beta_i]. \end{aligned}$$

The term $\sum_{i \in \alpha} T$ is equal to $|\alpha|T$. Since \mathcal{A} is laminar, for each $i, i' \in \alpha$ either $\beta_i = \beta_{i'}$ or $\beta_i \cap \beta_{i'} = \emptyset$, and then

$$\begin{aligned} \sum_{i \in \alpha} \text{TOT-LOAD}[i, \beta_i] &= \sum_{\substack{\beta \subset \alpha \\ \beta \text{ is maximal}}} \sum_{i \in \beta} \text{TOT-LOAD}[i, \beta] \\ &= \sum_{\substack{\beta \subset \alpha \\ \beta \text{ is maximal}}} \sum_{\gamma \subseteq \beta} \sum_{j=1}^n p_{\gamma j} x_{\gamma j} \\ &= \sum_{\gamma \subset \alpha} \sum_{j=1}^n p_{\gamma j} x_{\gamma j}, \end{aligned}$$

where the last two equalities follow from (3) and from the fact that \mathcal{A} is laminar, respectively. Therefore,

$$\begin{aligned} V &= \sum_{j=1}^n p_{\alpha j} x_{\alpha j} - |\alpha|T + \sum_{\gamma \subset \alpha} \sum_{j=1}^n p_{\gamma j} x_{\gamma j} \\ &= \sum_{\gamma \subset \alpha} \sum_{j=1}^n p_{\gamma j} x_{\gamma j} - |\alpha|T. \end{aligned}$$

Since $V > 0$, then $\sum_{\gamma \subset \alpha} \sum_{j=1}^n p_{\gamma j} x_{\gamma j} > |\alpha|T$, a contradiction to constraint (2b). \square

Algorithm 2 guarantees that for any set β there exists at most one machine $i \in \beta$ whose load is due to jobs assigned to α and to β , where α is some set such that $\beta \subset \alpha$. This is proven in the next lemma and will be exploited by the second phase of the algorithm.

ALGORITHM 3: Second phase (top-down job scheduling)

```
1 MARK[ $\alpha$ ]  $\leftarrow$  false for each  $\alpha \in \mathcal{A}$ ;  
2 while  $\exists \beta \in \mathcal{A}$  such that  $\neg$ MARK[ $\beta$ ] do  
3   Let  $\beta$  such that  $\neg$ MARK[ $\beta$ ] and (MARK[ $\alpha$ ], for each  $\alpha$  such that  $\beta \subset \alpha$ );  
4   if  $\exists i \in \beta$  such that LOAD[ $i, \beta$ ]  $> 0$  and LOAD[ $i, \alpha$ ]  $> 0$ , for some set  $\alpha \in \mathcal{A}$   
   such that  $\beta \subset \alpha$  then  
5     Let  $\alpha$  be the minimal set such that  $\beta \subset \alpha$  and LOAD[ $i, \alpha$ ]  $> 0$ ;  
6      $t_\beta \leftarrow t_{i\alpha}$ ;  
7      $\ell \leftarrow i$ ;  
8   else  
9      $t_\beta \leftarrow 0$ ;  
10     $\ell \leftarrow \min \beta$ ;  
11  end  
12  foreach  $k \in \beta$  in any order starting from  $\ell$  do  
13    Assign LOAD[ $k, \beta$ ] units of time of jobs  $j$  such that  $x_{\beta j} = 1$  to machine  $k$ ,  
    in the interval  $[t_\beta, t_\beta + \text{LOAD}[k, \beta] \pmod{T}]$ ;  
14     $t_\beta \leftarrow t_\beta + \text{LOAD}[k, \beta] \pmod{T}$ ;  
15     $t_{k\beta} \leftarrow t_\beta$ ;  
16  end  
17  MARK[ $\beta$ ]  $\leftarrow$  true;  
18 end
```

Lemma 4.2. *For each set $\beta \in \mathcal{A}$ there exists at most one machine $i \in \beta$ such that, for some set $\alpha \in \mathcal{A}$ such that $\beta \subset \alpha$, it holds that $\text{LOAD}[i, \beta] > 0$ and $\text{LOAD}[i, \alpha] > 0$.*

Proof. By contradiction, let us assume that there exist two machines i and i' , $i < i'$, such that $\text{LOAD}[i, \beta] > 0$, $\text{LOAD}[i, \alpha] > 0$, $\text{LOAD}[i', \beta] > 0$, and $\text{LOAD}[i', \alpha'] > 0$, for some $\alpha, \alpha' \in \mathcal{A}$ such that $\beta \subset \alpha, \alpha'$. Let us consider line 10 of Algorithm 2 at the iteration related to set β and machine i and let γ be the maximal set $\gamma \subset \beta$ such that $i \in \gamma$, that is $\text{LOAD}[i, \beta] = \min\{V, T - \text{TOT-LOAD}[i, \gamma]\}$.

- If $\min\{V, T - \text{TOT-LOAD}[i, \gamma]\} = V$, then $\text{LOAD}[i, \beta] = V$ and at the end of the iteration the algorithm sets $V = 0$ (line 12). Therefore, for each machine $i'' \in \beta$, $i'' > i$, the instruction at line 10 sets $\text{LOAD}[i'', \beta] = V = 0$. Since $i' > i$, then $\text{LOAD}[i', \beta] = 0$, a contradiction to $\text{LOAD}[i', \beta] > 0$.
- If $\min\{V, T - \text{TOT-LOAD}[i, \gamma]\} = T - \text{TOT-LOAD}[i, \gamma]$, then $\text{LOAD}[i, \beta] = T - \text{TOT-LOAD}[i, \gamma]$ and the algorithm sets $\text{TOT-LOAD}[i, \beta] = \text{TOT-LOAD}[i, \gamma] + \text{LOAD}[i, \beta] = T$ at line 11. Therefore, for each α such that $\beta \subset \alpha$, $\text{TOT-LOAD}[i, \alpha] = \text{TOT-LOAD}[i, \beta] = T$ and $\text{LOAD}[i, \alpha] = \min\{V, 0\} = 0$, a contradiction to $\text{LOAD}[i, \alpha] > 0$. \square

The pseudo-code of the second phase is given in Algorithm 3. As in Algorithm 2, variable MARK[α] is used to determine whether a set α has been visited or not. In this case, the algorithm visits all the sets in \mathcal{A} in top-down order (see the while loop at lines 2–18). Variable $t_{i\alpha}$ stores the latest time instant in

which a job assigned to set α is scheduled on machine $i \in \alpha$. Let β be a maximal set that has not been visited yet (line 3). By Lemma 4.2, there exists at most one machine $i \in \beta$ such that $\text{LOAD}[i, \beta] > 0$ and $\text{LOAD}[i, \alpha] > 0$, for some set $\alpha \in \mathcal{A}$ such that $\beta \subset \alpha$. If such a machine exists, then let α be the minimal set satisfying the previous condition (line 5) and let ℓ be the unique machine where both sets have some load (line 7). We first schedule jobs assigned to β from time t_{i_α} on machine ℓ and then we proceed by scheduling the remaining volume on the empty machines in β as done for global jobs in Algorithm 1. In detail, we initialize t_β to t_{i_α} (line 6); for each machine $k \in \beta$, starting from ℓ , we assign $\text{LOAD}[k, \beta]$ units of time of jobs assigned to β to machine k , in the interval $[t_\beta, t_\beta + \text{LOAD}[k, \beta] \pmod{T}]$ (line 13); and we update t_β and t_{k_β} by adding $\text{LOAD}[k, \beta] \pmod{T}$ (lines 14–15). In the case that there is no machine in β with some loads due to two different sets, the only difference is that ℓ is chosen as the smallest machine in β and t_β is initialized to 0 (lines 9–10).

Lemma 4.3. *In the schedule produced by Algorithm 3, all jobs receive the required amount of processing time.*

Proof. We show that for each α all the jobs j such that $x_{\alpha j} = 1$ receive the required amount of processing time, i.e., $\sum_{j=1}^n p_{\alpha j} x_{\alpha j}$. For each $i \in \alpha$, Algorithm 3 assigns $\text{LOAD}[i, \alpha]$ units of time to machine i and therefore assigns $\sum_{i \in \alpha} \text{LOAD}[i, \alpha]$ overall time to jobs j such that $x_{\alpha j} = 1$. By Lemma 4.1.ii, $\sum_{i \in \alpha} \text{LOAD}[i, \alpha] = \sum_{j=1}^n p_{\alpha j} x_{\alpha j}$. \square

Lemma 4.4. *In the schedule produced by Algorithm 3, no job is scheduled in parallel with itself.*

Proof. The proof is similar to that of Lemma 3.2. Assume by contradiction that a job j such that $x_{\beta j} = 1$ is scheduled on two different machines $i, i' \in \beta$ during the same time interval $[t_1, t_2]$, $t_2 > t_1$, and assume w.l.o.g. that j is scheduled first on machine i and then on machine i' . Then, in all the iterations of the loop at lines 12–16 of Algorithm 3, from that of i to that of i' , only job j is scheduled and thus, since $t_2 > t_1$, $p_{\beta j} > T$, a contradiction. \square

Lemma 4.5. *In the schedule produced by Algorithm 3, no machine processes more than one job in the same time interval.*

Proof. By contradiction, let us consider the first iteration of the loop at lines 12–16 in which a machine i that is already scheduling a job j in some time interval is assigned another job j' in the same interval.

Let us assume that $x_{\beta j'} = 1$. Moreover, all the jobs j'' such that $x_{\gamma j''} = 1$ for any $\gamma \subset \beta$ are not yet scheduled in the considered iteration, therefore $x_{\alpha j} = 1$, for some α such that $\beta \subset \alpha$.

By Lemma 4.2, i is the only machine such that $\text{LOAD}[i, \beta] > 0$ and $\text{LOAD}[i, \alpha] > 0$, then, i is the machine ℓ selected at line 7 of Algorithm 3. Let $\alpha = \alpha_0 \supset \alpha_1 \supset \alpha_2 \dots \supset \alpha_L = \beta$ be all the sets in \mathcal{A} such that $\beta \subseteq \alpha_p \subseteq \alpha$. We recall that $t_{i_{\alpha_p}}$ is the last time instant in which the algorithm schedules a job assigned to

α_p on machine i and that the value of t_β after line 6 of the algorithm is executed, but before line 14 is executed, is $\bar{t}_\beta = t_{\alpha_{L-1}} = t_{i\alpha} + \sum_{l=1}^{L-1} \text{LOAD}[i, \alpha_l] \pmod{T}$. Let t_x be the first time instant in which the algorithm schedules a job assigned to α on machine i , that is, jobs assigned to α are scheduled either in the interval $[t_x, t_{i\alpha}]$, if $t_x < t_{i\alpha}$, or in the intervals $[t_x, T]$ and $[0, t_{i\alpha}]$, otherwise. In the former case, to have that machine i processes a job in α and a job in β in the same time interval, we must have that $\bar{t}_\beta + \text{LOAD}[i, \beta] > t_x + T$, that is $t_{i\alpha} + \sum_{l=1}^L \text{LOAD}[i, \alpha_l] > t_x + T$. Since $t_{i\alpha} - t_x = \text{LOAD}[i, \alpha]$, this implies that $\sum_{l=0}^L \text{LOAD}[i, \alpha_l] > T$, a contradiction to Lemma 4.1.i. In the latter case, we must have that $\bar{t}_\beta + \text{LOAD}[i, \beta] > t_x$, where $\bar{t}_\beta = t_{i\alpha} + \sum_{l=1}^{L-1} \text{LOAD}[i, \alpha_l]$, that is $t_{i\alpha} + \sum_{l=1}^L \text{LOAD}[i, \alpha_l] > t_x$. Since $\text{LOAD}[i, \alpha] = T - t_x + t_{i\alpha}$, we obtain again the contradiction $\sum_{l=0}^L \text{LOAD}[i, \alpha_l] > T$. \square

Theorem 4.6. *Given a feasible solution (\mathbf{x}, T) to (IP-2), Algorithms 2 and 3 produce a valid schedule in the interval $[0, T]$.*

Proof. The statement follows from Lemmas 4.3–4.5. \square

5. Rounding strategy for the ILP

To round the fractional relaxation of (IP-2), we first transform it into a decision form by applying a standard pruning technique [9]. It suffices to decide, for an arbitrary but fixed value of T , the feasibility of the following system:

$$\sum_{\alpha \in \mathcal{A}} x_{\alpha j} = 1 \quad \text{for } j \in J \quad (\text{IP-3})$$

$$\sum_{j \in J} \sum_{\beta \subseteq \alpha} p_{\beta j} x_{\beta j} \leq |\alpha|T \quad \text{for } \alpha \in \mathcal{A} \quad (4a)$$

$$x_{\alpha j} \in \{0, 1\} \quad \text{for } \alpha \in \mathcal{A}, j \in J, \quad (4b)$$

$$x_{\alpha j} = 0 \quad \text{for } (\alpha, j) \notin R, \quad (4c)$$

where $R = \{(\alpha, j) \in \mathcal{A} \times J : p_{\alpha j} \leq T\}$. We have eliminated constraints (2c) by observing that they are satisfied by a 0-1 solution if and only if $p_{\alpha j} \leq T$ whenever $x_{\alpha j} = 1$. Therefore, one can simply set to zero all variables $x_{\alpha j}$ such that $p_{\alpha j} > T$, i.e., the variables with indices not in R . The binary search process for the minimal T for which (IP-3) is feasible requires a number of iterations logarithmic in the range of T , and therefore multiplies the overall running time by only a polynomial factor.

Without loss of generality, we can assume that the family \mathcal{A} always contains the singleton machine sets $\{1\}, \{2\}, \dots, \{m\}$; if not, these sets can be added to \mathcal{A} by setting the processing time of a job $j \in J$ on machine $i \in M$ as the processing time of j on the (inclusion-wise) minimal set in \mathcal{A} that contains i . Before discussing how to round the fractional relaxation of (IP-3), we show that a feasible fractional solution can always be modified so that, for every $\alpha \in \mathcal{A}$, $x_{\alpha j} = 0$ unless α is a singleton set. This follows easily by repeated application

of the next lemma, which allows to “push down” the fractional weights towards the singleton sets of the laminar family.

Lemma 5.1. *Let $\eta \in \mathcal{A}$ be a non-singleton set. If \mathbf{x} is a feasible solution to the LP relaxation of (IP-3), then there exists another feasible solution \mathbf{x}' to the same LP relaxation such that $x'_{\eta j} = 0$ and $x'_{\alpha j} = x_{\alpha j}$ whenever $\alpha \not\subseteq \eta$.*

Proof. For any $\alpha \in \mathcal{A}$, we define the *slack* of α in \mathbf{x} to be

$$\text{slack}(\alpha, \mathbf{x}) := |\alpha|T - \sum_{j \in J} \sum_{\beta \subseteq \alpha} p_{\beta j} x_{\beta j}.$$

Note that the LP relaxation of (IP-3) can be written as

$$\sum_{\alpha \in \mathcal{A}} x_{\alpha j} = 1 \quad \text{for } j \in J \quad (5a)$$

$$\text{slack}(\alpha, \mathbf{x}) \geq 0 \quad \text{for } \alpha \in \mathcal{A} \quad (5b)$$

$$x_{\alpha j} \geq 0 \quad \text{for } \alpha \in \mathcal{A}, j \in J \quad (5c)$$

$$x_{\alpha j} = 0 \quad \text{for } (\alpha, j) \notin R. \quad (5d)$$

Without loss of generality, assume that $\eta = \beta_1 \cup \dots \cup \beta_q$ with $\beta_1, \dots, \beta_q \in \mathcal{A}$, $\beta_1, \dots, \beta_q \subset \eta$, the sets β_1, \dots, β_q being maximal and pairwise disjoint. Because η has nonnegative slack in \mathbf{x} , we have

$$\begin{aligned} \sum_{j \in J} \sum_{\gamma \subseteq \beta_1} p_{\gamma j} x_{\gamma j} + \dots + \sum_{j \in J} \sum_{\gamma \subseteq \beta_q} p_{\gamma j} x_{\gamma j} + \sum_{j \in J} p_{\eta j} x_{\eta j} \\ \leq |\beta_1|T + \dots + |\beta_q|T, \end{aligned}$$

which is equivalent to

$$\sum_{j \in J} p_{\eta j} x_{\eta j} \leq \text{slack}(\beta_1, \mathbf{x}) + \dots + \text{slack}(\beta_q, \mathbf{x}). \quad (6)$$

We now define a new solution \mathbf{x}' by setting $x'_{\eta j} = 0$, $x'_{\alpha j} = x_{\alpha j}$ for $\alpha \neq \beta_1, \dots, \beta_q, \eta$, and

$$x'_{\beta j} = x_{\beta j} + \frac{\text{slack}(\beta, \mathbf{x})}{\text{slack}(\beta_1, \mathbf{x}) + \dots + \text{slack}(\beta_q, \mathbf{x})} \cdot x_{\eta j} \quad (7)$$

for $\beta = \beta_1, \dots, \beta_q$. We claim that \mathbf{x}' is valid for the LP. To see that (5a) is satisfied by the new solution, note that

$$\begin{aligned} \sum_{\alpha \in \mathcal{A}} x'_{\alpha j} &= \sum_{\substack{\alpha \in \mathcal{A} \\ \alpha \neq \eta}} x_{\alpha j} + \sum_{i=1}^q \frac{\text{slack}(\beta_i, \mathbf{x})}{\text{slack}(\beta_1, \mathbf{x}) + \dots + \text{slack}(\beta_q, \mathbf{x})} \cdot x_{\eta j} \\ &= \sum_{\alpha \in \mathcal{A}} x_{\alpha j} = 1. \end{aligned}$$

To see that (5b) is satisfied, it suffices to show that the new slack of β_1, \dots, β_q is nonnegative, since the slack of any other set does not decrease. Consider, say, β_i . By summing (7) across jobs,

$$\begin{aligned} \sum_{j \in J} \sum_{\gamma \subseteq \beta_i} p_{\gamma j} x'_{\gamma j} &= \sum_{j \in J} \sum_{\gamma \subseteq \beta_i} p_{\gamma j} x_{\gamma j} + \frac{\text{slack}(\beta_i, \mathbf{x})}{\text{slack}(\beta_1, \mathbf{x}) + \dots + \text{slack}(\beta_q, \mathbf{x})} \sum_{j \in J} p_{\beta_i, j} x_{\eta j} \\ &\leq \sum_{j \in J} \sum_{\gamma \subseteq \beta_i} p_{\gamma j} x_{\gamma j} + \frac{\text{slack}(\beta_i, \mathbf{x})}{\text{slack}(\beta_1, \mathbf{x}) + \dots + \text{slack}(\beta_q, \mathbf{x})} \sum_{j \in J} p_{\eta j} x_{\eta j} \\ &\leq \sum_{j \in J} \sum_{\gamma \subseteq \beta_i} p_{\gamma j} x_{\gamma j} + \text{slack}(\beta_i, \mathbf{x}), \end{aligned}$$

where for the first inequality we used the monotonicity of the processing times, and for the second inequality we used (6). Therefore,

$$\begin{aligned} \text{slack}(\beta_i, \mathbf{x}') &= \text{slack}(\beta_i, \mathbf{x}) - \sum_{j \in J} \sum_{\gamma \subseteq \beta_i} p_{\gamma j} x'_{\gamma j} + \sum_{j \in J} \sum_{\gamma \subseteq \beta_i} p_{\gamma j} x_{\gamma j} \\ &\geq \text{slack}(\beta_i, \mathbf{x}) - \text{slack}(\beta_i, \mathbf{x}) \geq 0. \end{aligned}$$

□

Equipped with the above lemma, we can proceed to prove an upper bound on the approximability of the hierarchical scheduling problem.

Theorem 5.2. *The hierarchical scheduling problem admits a polynomial-time 2-approximation algorithm.*

Proof. Consider an instance $I = (J, M, \mathcal{A}, p)$ of the hierarchical scheduling problem.

Let T^* be the minimum value of T for which the LP relaxation of (IP-3) is feasible. Clearly, $T^* \leq \text{opt}(I)$ where $\text{opt}(I)$ is the optimal makespan of the hierarchical scheduling instance I . By applying repeatedly Lemma 5.1, we can ensure that there exists a feasible fractional solution \mathbf{x} with makespan T^* and such that $x_{\alpha j} > 0$ only for α such that $|\alpha| = 1$. Because of this, observe that \mathbf{x} can also be seen as a fractional solution to an unrelated machines scheduling instance with makespan T^* : the instance $I_u = (J, M, p')$ obtained by defining $p'_{ij} := p_{\{i\}j}$.

The idea is now to invoke an existing LP-based algorithm for the unrelated machine scheduling problem and run it on I_u . The classic rounding algorithm by Lenstra, Shmoys and Tardos [15] constructs an integral assignment $\bar{\mathbf{x}}$ for I_u with makespan $T_u \leq 2T^*$. The assignment $\bar{\mathbf{x}}$, extended with 0 values on all sets with $|\alpha| > 1$, is also valid for (IP-3) if we take $T = T_u \leq 2T^* \leq 2 \cdot \text{opt}(I)$. Therefore, such an extended assignment yields a 2-approximate solution for the hierarchical scheduling problem. □

We note that the reduction used in the above proof, from fractional hierarchical scheduling to fractional unrelated machines, is *not* valid for the original

(integral) formulations of the two problems: indeed, in Example 2.1, the original instance I of the semi-partitioned problem has an optimal makespan of 2, while the unrelated machine instance I_u has an optimal makespan of 3. In general, the gap between the makespan of I_u and the makespan of I can be arbitrarily close to a factor 2, as the next example shows.

Example 5.1. (See also Figure 3.) Consider a semi-partitioned instance I with n jobs and $m := n - 1$ machines. Recall that we use machine index 0 to denote global processing. Job j , $j = 1, \dots, n - 1$, has $p_{ij} = n - 2$ if $i = j$, and $p_{ij} = \infty$ otherwise. Job n has $p_{ij} = n - 1$ for each $i = 0, 1, \dots, n - 1$. An optimal solution has makespan $\text{opt}(I) = n - 1$: assign job j , $j = 1, \dots, n - 1$ to machine j and assign job n globally; schedule each job j , $j = 1, \dots, n - 1$, on machine j in time intervals $[0, j - 1)$ and $[j, n - 1)$; schedule job n on machine i , $i = 1, \dots, n - 1$ during $[i - 1, i)$. On the other hand, in the corresponding unrelated machine instance I_u , jobs cannot be migrated and therefore the minimum value of the makespan is $2n - 3$.

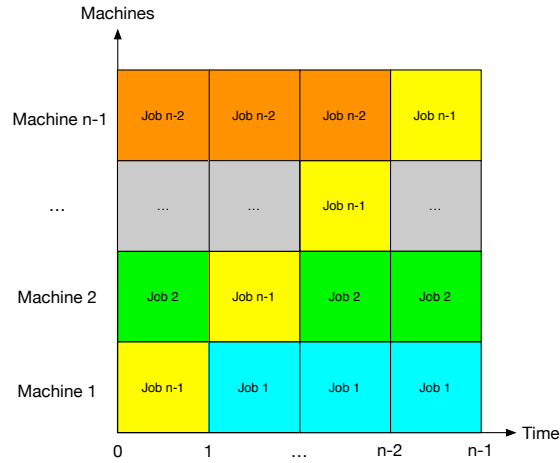


Figure 3: Optimal schedule for the instance of Example 5.1.

6. Memory constraints

The basic model as described in the previous sections focuses on makespan minimization, without additional constraints. However, machines often also have limited *memory capacity*. In this section we show how, in the hierarchical case, our model can be extended to incorporate memory capacities and discuss how to obtain efficient algorithms with a guaranteed bicriteria approximation ratio. We consider two distinct extensions, which we call Model 1 and Model

2. In the first model, each machine has a separate memory capacity. In the second model, each cluster of the hierarchical architecture has a certain memory capacity, which is shared among all machines of the cluster.

Model 1. We assume that each machine $i \in M$ has some memory budget $B_i \in \mathbb{Z}_+$ and that each job $j \in J$ requires memory space $s_{ij} \in \mathbb{Z}_+$ when run on machine i . We require the jobs assigned to sets that include machine i to fit the memory bound B_i ; i.e., if j is assigned to a set of machines α , then its space requirement is counted towards each machine in α . Thus, we revise ILP (IP-3) by adding the capacity constraints

$$\sum_{j \in J} \left(s_{ij} \cdot \sum_{\alpha \in \mathcal{A}: i \in \alpha} x_{\alpha j} \right) \leq B_i \quad \text{for each } i \in M. \quad (8)$$

To round the fractional relaxation of the revised ILP, we apply the *iterative rounding* approach [11, 13, 18], allowing us to prove the following theorem.

Theorem 6.1. *Whenever ILP (IP-3) has a solution satisfying constraints (8), it is possible to construct, in polynomial time, a valid schedule with makespan at most $3T$ such that*

$$\sum_{j \in J} \left(s_{ij} \cdot \sum_{\alpha \in \mathcal{A}: i \in \alpha} x_{\alpha j} \right) \leq 3 \cdot B_i \quad \text{for each } i \in M. \quad (9)$$

where \mathbf{x} represents the schedule's assignment.

Proof. Consider a feasible solution \mathbf{y} to the fractional relaxation of the revised ILP (IP-3), including the capacity constraints (8). By applying a similar reasoning as in the proof of Theorem 5.2, we can obtain a feasible fractional vector \mathbf{z} that satisfies the constraints of ILP (IP-3) and such that $z_{\alpha j} > 0$ only if $|\alpha| = 1$. Moreover, note that

$$z_{ij} \leq \sum_{\alpha \in \mathcal{A}: i \in \alpha} y_{\alpha j} \quad \text{for } i \in M, j \in J \quad (10)$$

because the fractional weight in \mathbf{z} associated to a pair (i, j) can only be due to the fractional weight in \mathbf{y} of pairs (α, j) such that $i \in \alpha$. Therefore,

$$\sum_{j \in J} s_{ij} z_{ij} \leq \sum_{j \in J} s_{ij} \left(\sum_{\alpha \in \mathcal{A}: i \in \alpha} y_{\alpha j} \right) \leq B_i, \quad (11)$$

where the second inequality follows by the feasibility of \mathbf{y} . This means that the

fractional solution \mathbf{z} is feasible for the following LP:

$$\sum_{i \in M} z_{ij} = 1 \quad \text{for } j \in J \quad (12a)$$

$$\sum_{j \in J} p_{ij} z_{ij} \leq T \quad \text{for } i \in M \quad (12b)$$

$$\sum_{j \in J} s_{ij} z_{ij} \leq B_i \quad \text{for } i \in M \quad (12c)$$

$$z_{ij} \geq 0 \quad \text{for } i \in M, j \in J, \quad (12d)$$

$$z_{ij} = 0 \quad \text{for } (i, j) \notin R, \quad (12e)$$

where $R = \{(i, j) \in M \times J : p_{ij} \leq T \wedge s_{ij} \leq b_i\}$. This linear program can be rounded using an existing iterated rounding technique, described in [18]. The algorithm of [18] constructs an integral assignment \mathbf{x} satisfying

$$\sum_{j \in J} p_{ij} x_{ij} \leq T + 2 \max_{j \in J: (i,j) \in R} p_{ij} \quad \text{for } i \in M \quad (13a)$$

$$\sum_{j \in J} s_{ij} x_{ij} \leq B_i + 2 \max_{j \in J: (i,j) \in R} s_{ij} \quad \text{for } i \in M, \quad (13b)$$

where the factor of 2 is due to the fact that each variable z_{ij} appears with at most two nonzero coefficients in the packing constraints (12b)-(12c) (one is p_{ij} , the other s_{ij}). Since $\max_{j \in J: (i,j) \in R} p_{ij} \leq T$ and $\max_{j \in J: (i,j) \in R} s_{ij} \leq B_i$, this proves the claim. \square

Model 2. Consider the forest associated to the laminar family \mathcal{A} , that is, the forest having as nodes the sets in \mathcal{A} , and such that α is a parent of β if and only if $\beta \subset \alpha$ and there is no $\gamma \in \mathcal{A}$ such that $\beta \subset \gamma \subset \alpha$. We can assume that this forest is, in fact, a tree: if it is not, we can add to \mathcal{A} the set M containing all machines, and set a very high value of p_{Mj} for each job j (i.e. $p_{Mj} = \infty$) to ensure that the set of feasible solutions is not affected.

Let $\mu : \mathcal{A} \rightarrow \mathbb{Q}_+$. In this model, we assume that job j requires space $s_j \leq 1$ (with $s_j \in \mathbb{Q}_+$), that each node α except the root has memory capacity $\mu(\alpha) \geq 1$, and that the root has unbounded capacity. Thus, μ is a function capturing how the memory hierarchy scales. We require the jobs assigned to a set (node of the tree) to fit the memory capacity of that set. Thus, we revise ILP (IP-3) by adding the capacity constraints

$$\sum_{j \in J} s_j x_{\alpha j} \leq \mu(\alpha) \quad \text{for each } \alpha \in \mathcal{A} \setminus \{M\}. \quad (14)$$

Call (IP-4) the revised ILP obtained in this way.

The additional constraints affect the applicability of Theorem 5.2, since one cannot use Lemma 5.1 to “push down” the values of the fractional variables towards the leaves of the laminar family. Moreover, the approximability bounds ensured by known rounding techniques (for example, [12, 18]) are not suitable in

this case. Indeed, one cannot apply the result of Karp et al. [12] because it does not ensure that the resulting integral vector satisfies the assignment constraints exactly; while the guarantee of Marchetti-Spaccamela et al. [18] yields a large approximation factor, equal to $1 + k$, where k is the number of levels of the hierarchy.

We improve the analysis of the iterative rounding scheme from Marchetti-Spaccamela et al. [18]. We show, in Lemma 6.2 below, that the scheme satisfies the assignment constraints exactly and yields a bound in terms of the *sum* of the column's entries of the (normalized) coefficient matrix; when applied to (IP-4), this guarantees $O(\log k)$ approximation of the packing constraints. Such a rounding scheme applies to general assignment and packing constraints, and thus may find applications beyond the hierarchical scheduling problem.

Lemma 6.2. *Let I, J be nonempty finite sets, and $R \subseteq I \times J$. Consider a linear program of the form:*

$$\min \sum_{(i,j) \in R} c_{ij} z_{ij} \quad (\text{LP})$$

$$\sum_{i:(i,j) \in R} z_{ij} = 1 \quad \forall j \in J \quad (15a)$$

$$\sum_{q=(i,j) \in R} a_{lq} z_{ij} \leq b_l \quad l = 1, \dots, \theta \quad (15b)$$

$$0 \leq z_{ij} \leq 1 \quad \forall (i,j) \in R, \quad (15c)$$

where z_{ij} are variables, $\theta \in \mathbb{N}$, and $a_{lq} \geq 0$, $b_l > 0$, $c_{ij} \geq 0$ for all $l = 1, \dots, \theta$, $q = (i,j) \in R$. Assume that the LP has a feasible solution \mathbf{z}^0 and that the bound $\sum_{l=1}^{\theta} a_{lq}/b_l \leq \rho$ holds for each $q \in R$. Then there are values \bar{z}_{ij} such that

$$\sum_{(i,j) \in R} c_{ij} \bar{z}_{ij} \leq \sum_{(i,j) \in R} c_{ij} z_{ij}^0 \quad (16a)$$

$$\sum_{i:(i,j) \in R} \bar{z}_{ij} = 1 \quad \forall j \in J \quad (16b)$$

$$\sum_{q=(i,j) \in R} a_{lq} \bar{z}_{ij} \leq (1 + \rho) b_l \quad l = 1, \dots, \theta \quad (16c)$$

$$\bar{z}_{ij} \in \{0, 1\} \quad \forall (i,j) \in R. \quad (16d)$$

Proof. Before formally proving the Lemma we present the iterative rounding procedure detailed in Algorithm 4. As in Marchetti-Spaccamela et al. [18], the idea of the algorithm is to compute, at iteration $h = 0, 1, 2, \dots$, an optimal extreme-point solution \mathbf{z}^h of a linear program LP^h , with LP^0 being the initial program (LP) assumed in the Lemma. Each subsequent LP is obtained by either freezing some variables at their integer value in the current LP solution (and updating the corresponding right-hand side coefficients – Lines 4–5), or, if

ALGORITHM 4: Iterative ILP rounding

```
1 while LP has  $s \geq 1$  variables do
2   Solve LP to find extreme-point solution  $\mathbf{z}^h$ ;
3   if  $\mathbf{z}^h$  has at least one integral entry then
4     fix all integral variables at their value in  $\mathbf{z}^h$ , remove them from LP and
5     update right-hand side coefficients;
6     remove from LP all constraints that no longer involve any variable;
7   else
8     find a constraint index  $l \in \{1, \dots, \theta\}$  such that
9      $\sum_{q=(i,j) \in R} a_{lq}(1 - z_{ij}^h) \leq \rho \cdot b_l$ , where  $S = \{0, 1\}^s$ ;
10    remove from LP the constraint corresponding to  $l$ ;
11  end
12 end
```

no variable is integral, by discarding a packing constraint $l \in \{1, \dots, \theta\}$ among those in (15b) that satisfy

$$\sum_{q=(i,j) \in R} a_{lq}(1 - z_{ij}^h) \leq \rho \cdot b_l,$$

with $S = \{0, 1\}^s$. Such a packing constraint can be “safely dropped”, in the sense that even if the partial integer solution is completed by setting all remaining variables to their “worst” possible value, which is 1, the constraint will eventually be violated by at most ρ times its original right-hand side b_l (Lines 7–8). Note that the assignment constraints (15a) are never dropped, and that if z_{ij}^h is fixed at value 1, all remaining variables $z_{i'j}^h$ with $i' \neq i$ are fixed at value 0, due to the structure of the assignment constraints.

A crucial point, of course, is to guarantee that whenever the solution of LP^h has only fractional entries (and thus the **else** branch is taken in Line 6 of the algorithm), then there always exists some constraint of LP^h that can be safely dropped, i.e., the appropriate index l can always be found in Line 7 of the algorithm, ensuring progress of the iterative procedure. This fact is proven in the following auxiliary lemma.

Lemma 6.3. *Let LP^h be the linear program that is solved in iteration h of the rounding procedure, having s variables and r constraints. Let \mathbf{z}^h be an extreme-point solution to this LP. If $s > r$, then \mathbf{z}^h has at least one integral entry. If $s \leq r$, then there is some $l \in \{1, \dots, \theta\}$ such that $\sum_{q=(i,j) \in R} a_{lq}(1 - z_{ij}^h) \leq \rho \cdot b_l$, where S is the integer solution space for all remaining variables, i.e., $S = \{0, 1\}^s$.*

Proof. Assume that the subproblem in iteration h (described by LP^h) is defined as $\mathbf{Az} \leq \mathbf{b}$, where $\mathbf{A} \in \mathbb{R}^{r \times s}$. Note that $r = r_a + r_b$, where r_a is the number of constraints of type (15a) and r_b is the number of constraints of type (15b). We can ignore the constraints $0 \leq z_{ij} \leq 1$ because if any of them is active,

the iteration immediately completes with one variable being fixed to 0 or 1 and removed from the problem.

First assume that $s > r$. Then the null space of \mathbf{A} is nontrivial, so let \mathbf{z}_0 be a nonzero vector in the null space of \mathbf{A} . Since \mathbf{z}^h is an extreme-point solution to LP^h , it cannot be expressed as the convex combination of two (or more) solutions to LP^h . If \mathbf{z}^h does not have any integral entry, then we can find a value $\delta > 0$ such that $\mathbf{z}^h + \delta\mathbf{z}_0$ and $\mathbf{z}^h - \delta\mathbf{z}_0$ are both solutions to LP^h and, in particular, \mathbf{z}^h is a convex combination of these two solutions. Therefore \mathbf{z}^h must have at least one integral entry.

Now assume that $s \leq r$. For this case, we show that there always exists a constraint l of type (15b) such that $\{(\mathbf{A}\mathbf{1})_l - (\mathbf{A}\mathbf{z})_l\} \leq \rho \cdot b_l$. We show the statement by contradiction. Assume that the statement is not true; that is, for each constraint l of type (15b) it holds that

$$(\mathbf{A}\mathbf{1})_l - (\mathbf{A}\mathbf{z})_l > \rho \cdot b_l. \quad (17)$$

In each previous round, if variables were removed from the program, also constraints that had become redundant, were removed. Therefore, for all variables present in the linear program in this round, the corresponding constraint of type (15a) is also still present in the linear program (and this constraint is not present if all of its variables have been removed from the program). It follows that

$$\sum_{(i,j) \in R^h} z_{ij} = r_a, \quad (18)$$

where R^h is the index set of the variables in LP^h . Define Θ^h as the set of constraints of type (15b) present in the current linear program LP^h . Then,

$$\begin{aligned} \rho(r - r_a) &= \rho r_b \\ &< \sum_{l \in \Theta^h} \frac{1}{b_l} ((\mathbf{A}\mathbf{1})_l - (\mathbf{A}\mathbf{z})_l) \\ &= \sum_{l \in \Theta^h} \frac{1}{b_l} \sum_{q \in R^h} a_{lq} (1 - z_q) \\ &= \sum_{q \in R^h} (1 - z_q) \sum_{l \in \Theta^h} \frac{a_{lq}}{b_l} \\ &\leq \sum_{q \in R^h} \rho (1 - z_q) \\ &= \rho s - \sum_{q \in R^h} \rho z_q \\ &= \rho(s - r_a). \end{aligned}$$

The second inequality follows by the assumption on the normalized column sums of \mathbf{A} .

The chain of inequalities implies that $\rho(r - r_a) < \rho(s - r_a)$, that is, $r < s$, which is a contradiction to being in the case that $s \leq r$. Hence, we conclude that

if $s \leq r$, there must be a constraint l of type (15b), for which $\{(\mathbf{A}\mathbf{1})_l - (\mathbf{A}\mathbf{z})_l\} \leq \rho \cdot b_l$. \square

We now complete the proof of Lemma 6.2 showing that the obtained solution verifies the claim. Lemma 6.3 guarantees that when the extreme-point solution \mathbf{z}^h has no integral entries, one of the constraints (15b) can be dropped without violating the corresponding guarantee in (16c) in subsequent steps: even if all remaining variables are set to 1, (16c) will be satisfied for that value of l . Therefore, Algorithm 4 always finds either an integral variable – which is fixed and removed – or an appropriate constraint that is discarded (i.e., Line 7 always finds an appropriate l).

By induction, each LP^h is feasible: LP^0 has feasible solution \mathbf{z}^0 by assumption, and each subsequent LP is obtained by discarding constraints or by fixing some variables at their current LP value; both operations preserve feasibility, and do not increase the cost of the optimal solution; indeed, \mathbf{z}^h induces a solution to LP^{h+1} with the same objective value as \mathbf{z}^h in LP^h . Thus, the final vector $\bar{\mathbf{z}}$ fixed by the iterative rounding process has 0/1 components, satisfies the assignment constraints (i.e., (16b) holds), and has a cost bounded by the cost of the initial feasible solution \mathbf{z}^0 (i.e., (16a) holds). In general, $\bar{\mathbf{z}}$ may violate some of the packing constraints, but by the choice of the dropped constraints in Algorithm 4, no packing constraint l on $\bar{\mathbf{z}}$ will be violated by more than an additional factor $\rho \cdot b_l$, i.e., (16c) holds. This concludes the proof of Lemma 6.2. \square

Equipped with the rounding lemma, we can proceed to prove our result for Model 2.

Theorem 6.4. *Let k be the number of levels of the laminar family \mathcal{A} and let H_k be the k th harmonic number. Whenever ILP (IP-4) has a solution, it is possible to construct, in polynomial time, a valid schedule with makespan at most $\sigma \cdot T$ such that*

$$\sum_{j \in J} s_j x_{\alpha j} \leq \sigma \cdot \mu(\alpha) \text{ for each } \alpha \in \mathcal{A} \setminus \{M\}, \quad (19)$$

where \mathbf{x} represents the schedule's assignment and $\sigma = 2 + H_k$. When $k = 2$, the same holds with $\sigma = 3 + 1/m$.

Proof. We observe that (IP-4) is in a form suitable for applying Lemma 6.2. In particular, we take I in Lemma 6.2 to be the admissible sets family \mathcal{A} , and we use the generic packing constraints (16c) to encode constraints (4a) and (14). Indeed, constraints (4a) can be encoded as $|\mathcal{A}|$ constraints with coefficients of the form

$$a_{\alpha,(\beta,j)} := \begin{cases} p_{\beta j} & \text{if } \beta \subseteq \alpha, \\ 0 & \text{otherwise,} \end{cases}, \quad b_{\alpha} := |\alpha|T, \quad (\alpha \in \mathcal{A}),$$

while constraints (14) can be encoded as $|\mathcal{A}| - 1$ constraints with coefficients of the form

$$a_{\alpha,(\beta,j)} := \begin{cases} s_j & \text{if } \beta = \alpha, \\ 0 & \text{otherwise,} \end{cases}, b_\alpha := \mu(\alpha), (\alpha \in \mathcal{A} \setminus \{M\}).$$

The cost coefficients c_{ij} in (LP) can be set to zero and (LP) becomes the linear relaxation of (IP-4). By our hypothesis, such LP relaxation must be feasible, and the hypothesis of Lemma 6.2 is satisfied. In particular, we can take \mathbf{z}^0 to be an optimal solution of the LP relaxation. Note that the LP relaxation can be solved in polynomial time; in particular, since \mathcal{A} is laminar, $|\mathcal{A}| \leq 2m$ (see, e.g., [24, Theorem 3.5]) and the number of constraints in (IP-4) is polynomial in n and m .

To choose an appropriate value of ρ in Lemma 6.2, note that $(\beta, j) \in R$ only when $p_{\beta j} \leq T$ by construction, so if $q = (\beta, j)$,

$$\begin{aligned} \sum_{l=1}^{\theta} \frac{a_{lq}}{b_l} &= \sum_{\alpha \in \mathcal{A}: \beta \subseteq \alpha} \frac{p_{\beta j}}{|\alpha|T} + \frac{s_j}{\mu(\beta)} \leq \sum_{\alpha \in \mathcal{A}: \beta \subseteq \alpha} \frac{1}{|\alpha|} + 1 \\ &\leq 1 + \sum_{1 \leq i \leq k} \frac{1}{i} = 1 + H_k, \end{aligned}$$

where for the first inequality we also used $s_j \leq 1 \leq \mu(\beta)$, and for the second the fact that the laminar family \mathcal{A} has k levels and the fact that all sets in \mathcal{A} are distinct and nonempty. Thus, we can apply Lemma 6.2 with $\rho = 1 + H_k$.

In the semi-partitioned case (i.e., when \mathcal{A} has $k = 2$ levels), the summation $\sum_{l=1}^{\theta} a_{lq}/b_l$ involves at most three nonzero terms. The first term is due to the local scheduling constraints and has the form p_{ij}/T , which is at most 1. The second term is due to the global scheduling constraint and has the form p_{ij}/mT , which is at most $1/m$. The third term is due to the memory constraints and has the form $s_j/\mu(\beta)$, which is at most 1. Therefore, $\rho = 2 + 1/m$ is sufficient.

Finally, the integer solution \mathbf{x} to (16), which can be found by applying Lemma 6.2, can be used to construct a schedule with makespan at most $(1 + \rho)T$ and satisfying (19) by feeding \mathbf{x} to the algorithms presented in Sections 3 and 4. \square

Acknowledgment

We thank an anonymous reviewer for pointing out references [7, 17] and their connection with the non-hierarchical setting.

References

- [1] Sanjoy Baruah. 2015. Federated Scheduling of Sporadic DAG Task Systems. In *Proc. 2015 IEEE International Parallel and Distributed Processing Symposium*. 179–186. <https://doi.org/10.1109/IPDPS.2015.33>

- [2] Andrea Bastoni, Björn B. Brandenburg, and James H. Anderson. 2010. An Empirical Comparison of Global, Partitioned, and Clustered Multiprocessor EDF Schedulers. In *Proc. of the 31st IEEE Real-Time Systems Symposium*. IEEE, 14–24. <https://doi.org/10.1109/RTSS.2010.23>
- [3] Andrea Bastoni, Björn B. Brandenburg, and James H. Anderson. 2011. Is Semi-Partitioned Scheduling Practical?. In *Proc. of the 23rd Euromicro Conf. on Real-Time Systems*. IEEE, 125–135. <https://doi.org/10.1109/ECRTS.2011.20>
- [4] Guy E. Blelloch, Jeremy T. Fineman, Phillip B. Gibbons, and Harsha Vardhan Simhadri. 2011. Scheduling irregular parallel computations on hierarchical caches. In *Proc. of the 23rd ACM Symposium on Parallelism in Algorithms and Architectures*. ACM, 355–366. <https://doi.org/10.1145/1989493.1989553>
- [5] Vincenzo Bonifaci, Björn B. Brandenburg, Gianlorenzo D’Angelo, and Alberto Marchetti-Spaccamela. 2016. Multiprocessor Real-Time Scheduling with Hierarchical Processor Affinities. In *Proc. 28th Euromicro Conf. on Real-Time Systems*. 237–247. <https://doi.org/10.1109/ECRTS.2016.24>
- [6] Marin Bougeret, Pierre-Francois Dutot, Denis Trystram, Klaus Jansen, and Christina Robenek. 2015. Improved approximation algorithms for scheduling parallel jobs on identical clusters. *Theoretical Computer Science* 600 (2015), 70 – 85. <https://doi.org/10.1016/j.tcs.2015.07.003>
- [7] José R. Correa, Martin Skutella, and José Verschae. 2012. The Power of Preemption on Unrelated Machines and Applications to Scheduling Orders. *Math. Oper. Res.* 37, 2 (2012), 379–398. <https://doi.org/10.1287/moor.1110.0520>
- [8] Celia A. Glass and Hans Kellerer. 2007. Parallel machine scheduling with job assignment restrictions. *Naval Research Logistics* 54, 3 (2007), 250–257. <https://doi.org/10.1002/nav.20202>
- [9] Dorit S. Hochbaum and David B. Shmoys. 1987. Using dual approximation algorithms for scheduling problems: theoretical and practical results. *J. ACM* 34, 1 (1987), 144–162. <https://doi.org/10.1145/7531.7535>
- [10] Jing-Jang Hwang, Yuan-Chieh Chow, Frank D. Anger, and Chung-Yee Lee. 1989. Scheduling Precedence Graphs in Systems with Interprocessor Communication Times. *SIAM J. Comput.* 18, 2 (1989), 244–257. <https://doi.org/10.1137/0218016>
- [11] Kamal Jain. 2001. A Factor 2 Approximation Algorithm for the Generalized Steiner Network Problem. *Combinatorica* 21, 1 (2001), 39–60. <https://doi.org/10.1007/s004930170004>

- [12] Richard M. Karp, Frank Thomson Leighton, Ronald L. Rivest, Clark D. Thompson, Umesh V. Vazirani, and Vijay V. Vazirani. 1987. Global Wire Routing in Two-Dimensional Arrays. *Algorithmica* 2 (1987), 113–129. <https://doi.org/10.1007/BF01840353>
- [13] L.C. Lau, R. Ravi, and M. Singh. 2011. *Iterative Methods in Combinatorial Optimization*. Cambridge University Press.
- [14] Eugene L. Lawler and Jacques Labetoulle. 1978. On Preemptive Scheduling of Unrelated Parallel Processors by Linear Programming. *J. ACM* 25, 4 (1978), 612–619. <https://doi.org/10.1145/322092.322101>
- [15] Jan Karel Lenstra, David B. Shmoys, and Éva Tardos. 1990. Approximation Algorithms for Scheduling Unrelated Parallel Machines. *Math. Program.* 46 (1990), 259–271. <https://doi.org/10.1007/BF01585745>
- [16] Jing Li, Jian-Jia Chen, Kunal Agrawal, Chenyang Lu, Christopher D. Gill, and Abusayeed Saifullah. 2014. Analysis of Federated and Global Scheduling for Parallel Real-Time Tasks. In *Proc. 26th Euromicro Conf. on Real-Time Systems*. 85–96. <https://doi.org/10.1109/ECRTS.2014.23>
- [17] Jyh-Han Lin and Jeffrey Scott Vitter. 1992. ϵ -Approximations with Minimum Packing Constraint Violation (Extended Abstract). In *Proc. 24th Annual ACM Symposium on Theory of Computing*. ACM, 771–782.
- [18] Alberto Marchetti-Spaccamela, Cyriel Rutten, Suzanne van der Ster, and Andreas Wiese. 2015. Assigning sporadic tasks to unrelated machines. *Math. Program.* 152, 1-2 (2015), 247–274. <https://doi.org/10.1007/s10107-014-0786-9>
- [19] Evangelos P. Markatos and Thomas J. LeBlanc. 1994. Using Processor Affinity in Loop Scheduling on Shared-Memory Multiprocessors. *IEEE Trans. Parallel Distrib. Syst.* 5, 4 (1994), 379–400. <https://doi.org/10.1109/71.273046>
- [20] Robert McNaughton. 1959. Scheduling with Deadlines and Loss Functions. *Management Science* 6, 1 (1959), 1–12. <https://doi.org/10.1287/mnsc.6.1.1>
- [21] Gabriella Muratore, Ulrich M. Schwarz, and Gerhard J. Woeginger. 2010. Parallel machine scheduling with nested job assignment restrictions. *Operations Research Letters* 38, 1 (2010), 47 – 50. <https://doi.org/10.1016/j.orl.2009.09.010>
- [22] Rolf Rabenseifner, Georg Hager, and Gabriele Jost. 2009. Hybrid MPI/OpenMP Parallel Programming on Clusters of Multi-Core SMP Nodes. In *Proc. of the 17th Euromicro International Conf. on Parallel, Distributed and Network-Based Processing*. IEEE, 427–436. <https://doi.org/10.1109/PDP.2009.43>

- [23] James D. Salehi, James F. Kurose, and Donald F. Towsley. 1996. The effectiveness of affinity-based scheduling in multiprocessor network protocol processing (extended version). *IEEE/ACM Trans. Netw.* 4, 4 (1996), 516–530. <https://doi.org/10.1109/90.532862>
- [24] Alexander Schrijver. 2003. *Combinatorial Optimization – Polyhedra and Efficiency*. Springer.
- [25] Yanyong Zhang, Hubertus Franke, José E. Moreira, and Anand Sivasubramanian. 2003. An Integrated Approach to Parallel Scheduling Using Gang-Scheduling, Backfilling, and Migration. *IEEE Trans. Parallel Distrib. Syst.* 14, 3 (2003), 236–247. <https://doi.org/10.1109/TPDS.2003.1189582>
- [26] Sergey Zhuravlev, Juan Carlos Saez, Sergey Blagodurov, Alexandra Fedorova, and Manuel Prieto. 2012. Survey of scheduling techniques for addressing shared resources in multicore processors. *ACM Comput. Surv.* 45, 1 (2012), 4. <https://doi.org/10.1145/2379776.2379780>