

## APPROXIMATION ALGORITHMS FOR SCHEDULING UNRELATED PARALLEL MACHINES

Jan Karel LENSTRA

*Eindhoven University of Technology, Eindhoven, The Netherlands, and  
Centre for Mathematics and Computer Science, Amsterdam, The Netherlands*

David B. SHMOYS and Éva TARDOS

*Cornell University, Ithaca, NY, USA*

Received 1 October 1987

Revised manuscript 26 August 1988

We consider the following scheduling problem. There are  $m$  parallel machines and  $n$  independent jobs. Each job is to be assigned to one of the machines. The processing of job  $j$  on machine  $i$  requires time  $p_{ij}$ . The objective is to find a schedule that minimizes the makespan.

Our main result is a polynomial algorithm which constructs a schedule that is guaranteed to be no longer than twice the optimum. We also present a polynomial approximation scheme for the case that the number of machines is fixed. Both approximation results are corollaries of a theorem about the relationship of a class of integer programming problems and their linear programming relaxations. In particular, we give a polynomial method to round the fractional extreme points of the linear program to integral points that nearly satisfy the constraints.

In contrast to our main result, we prove that no polynomial algorithm can achieve a worst-case ratio less than  $\frac{3}{2}$  unless  $P = NP$ . We finally obtain a complexity classification for all special cases with a fixed number of processing times.

*Key words:* Scheduling, parallel machines, approximation algorithm, worst case analysis, linear programming, integer programming, rounding.

### 1. Introduction

Although the performance of approximation algorithms has been studied for over twenty years, very little is understood about the structural properties of a problem that permit good performance guarantees. In fact, there are practically no tools to distinguish those problems for which there does exist a polynomial algorithm for any performance bound, and those for which this is not the case. One problem area in which these questions have received much attention is that of scheduling and bin packing. We examine a scheduling problem for which all previously analyzed polynomial algorithms have particularly poor performance guarantees. We present

A preliminary version of this paper appeared in the *Proceedings of the 28th Annual IEEE Symposium on the Foundations of Computer Science* (Computer Society Press of the IEEE, Washington, D.C., 1987) pp. 217–224.

a polynomial algorithm that delivers a solution guaranteed to be within a factor of 2 of the optimum, and prove that this is nearly best possible, in the sense that no polynomial algorithm can guarantee a factor less than  $\frac{3}{2}$  unless  $P = NP$ . Our algorithm is based on a result concerning the relationship of certain integer programming problems and their linear relaxations that is of interest in its own right.

One of the most natural strategies to obtain good solutions to an integer linear program is to drop the integrality constraints, solve the resulting linear programming problem, and then round the solution to an integral solution. There are many difficulties with this approach. The rounded solution may be suboptimal, and it may even be impossible to round the solution to a feasible solution. For restricted classes of integer programs, however, the behavior might not be quite as bad. Certainly, if the extreme points of the linear programming relaxation are all integral, then the optimal solution is obtained without even rounding, as is the case, for example, for the bipartite matching, maximum flow, and minimum cost flow problems.

It is an interesting question to study those classes of integer programs for which the linear relaxations provide a good approximation, in that rounded solutions can be found that are nearly feasible or nearly optimal. Much work along these lines has been done for integer programs where the coefficients of the constraints are restricted to  $\{0, 1\}$  (Lovász, 1975; Chvátal, 1979; Bartholdi, Orlin and Ratliff, 1980; Bartholdi, 1981; Baum and Trotter, 1981; Marcotte, 1983; Aharoni, Erdős and Linial, 1985; Raghavan and Thompson, 1985; Raghavan, 1986). We present a rounding theorem of this sort for a natural class of integer programs with arbitrary coefficients.

The scheduling problem to be considered is as follows. There are  $n$  independent jobs that must be scheduled without preemption on a collection of  $m$  parallel machines. If job  $j$  is scheduled on machine  $i$ , the processing time required is  $p_{ij}$ , which we assume to be a positive integer. The total time used by machine  $i$  is the sum of the  $p_{ij}$  for the jobs that are assigned to machine  $i$ , and the makespan of a schedule is the maximum total time used by any machine. The objective is to find a schedule that minimizes the makespan. Graham, Lawler, Lenstra and Rinnooy Kan (1979) denote this problem by  $R||C_{\max}$ . Davis and Jaffe (1981) presented a list scheduling algorithm and proved that it delivers a schedule with makespan no more than  $2\sqrt{m}$  times the optimum. Until now, no polynomial algorithm with a better performance bound was known. We present a polynomial algorithm that guarantees a factor of 2.

Approximation algorithms for this problem and several of its special cases have been studied for over two decades. Much of this work has focused on the case where the machines are *identical*; that is,  $p_{hj} = p_{ij}$  for any job  $j$  and any two machines  $h, i$ . The area of worst-case analysis of approximation algorithms for NP-hard optimization problems can be traced to Graham (1966), who showed that for this special case with identical machines, a list scheduling algorithm always delivers a schedule with makespan no more than  $(2 - 1/m)$  times the optimum. We shall refer to an algorithm that is guaranteed to produce a solution of length no more than  $\rho$  times the optimum as a  $\rho$ -*approximation algorithm*. Note that we do not require such

an algorithm to be polynomial, although our primary focus will be on this subclass.

An important family of further restricted cases is obtained by considering a fixed number of identical machines. Graham (1969) showed that for any specified number  $m$  of machines, it is possible to obtain a polynomial  $(1 + \varepsilon)$ -approximation algorithm for any fixed  $\varepsilon > 0$ , but the running time depends exponentially on  $1/\varepsilon$  (and on  $m$ ). Such a family of algorithms is called a *polynomial approximation scheme*. This result was improved by Sahni (1976), who reduced the dependence of the running time on  $1/\varepsilon$  to a polynomial. Such a family of algorithms is called a *fully polynomial approximation scheme*.

If the number of machines is specified as part of the problem instance, results by Garey and Johnson (1975, 1978) imply that no fully polynomial approximation scheme can exist, even if the machines are identical, unless  $P = NP$ . However, Hochbaum and Shmoys (1987) presented a polynomial approximation scheme for the problem with identical machines.

A natural generalization of identical machines is the case of machines that run at different speeds but do so *uniformly*. Thus, for each machine  $i$  there is a speed factor  $s_i$ , and  $p_{ij} = p_j/s_i$  where  $p_j$  is the inherent processing requirement of job  $j$ . Results analogous to the case of identical machines have been obtained for uniform machines. Gonzalez, Ibarra and Sahni (1977) gave a polynomial 2-approximation algorithm. For any fixed number of machines, Horowitz and Sahni (1976) presented a fully polynomial approximation scheme, whereas Hochbaum and Shmoys (1988) gave a polynomial approximation scheme for the case that the number of machines is a part of the problem instance.

Given these strong results for special cases, there was no apparent reason to suspect that analogous results did not hold for the general setting of *unrelated* machines. In fact, Horowitz and Sahni (1976) also presented a fully polynomial approximation scheme for any fixed number of unrelated machines. However, for the case that the number of machines is specified as part of the problem instance, a polynomial approximation scheme is unlikely to exist. We prove that the existence of a polynomial  $(1 + \varepsilon)$ -approximation algorithm for any  $\varepsilon < \frac{1}{2}$  would imply that  $P = NP$ .

An interesting algorithm for this problem was presented by Potts (1985). It is a 2-approximation algorithm with running time bounded by  $m^{m-1}$  times a polynomial in the input size. At first glance, his result does not appear to be particularly interesting, since for fixed  $m$ , a fully polynomial approximation scheme was already known. However, that scheme not only requires time  $O(nm(nm/\varepsilon)^{m-1})$  but also space  $O((nm/\varepsilon)^{m-1})$ , while Potts' algorithm requires only polynomial space. Thus, from both a practical and a theoretical viewpoint, Potts' algorithm is a valuable contribution. It is based on extending the integral part of a linear programming solution by an enumerative process. We extend his work by proving that the fractional solution to the linear program can be rounded to a good integral approximation in polynomial time, thereby obviating the need for enumeration and removing the

exponential dependence on  $m$ . Our algorithm can deliver a schedule of length arbitrarily close to twice the optimum, even when there are only two machines. For the case  $m = 2$ , Potts (1985) presented a linear-time  $\frac{3}{2}$ -approximation algorithm, and Numata (1988) proposed a polynomial approximation scheme.

We also consider the problem with any fixed number of machines. We present a polynomial approximation scheme for this case, where the space required is bounded by a polynomial in the input,  $m$ , and  $\log(1/\varepsilon)$ .

Another natural way to restrict the problem is to consider instances where the number of different processing times is bounded. For example, if all processing times are equal, then the optimum schedule is computable in polynomial time. As a byproduct of our investigation, we obtain a complete characterization of the polynomially solvable special cases with a fixed number of processing times under the assumption that  $P \neq \text{NP}$ .

## 2. A rounding theorem

We first present the key tool for our approximation algorithms. Let  $J_i(t)$  denote the set of jobs that require no more than  $t$  time units on machine  $i$ , and let  $M_j(t)$  denote the set of machines that can process job  $j$  in no more than  $t$  time units. Consider a decision version of our scheduling problem, where for each machine  $i$  there is a deadline  $d_i$  and where we are further constrained to schedule jobs so that each uses processing time at most  $t$ ; we wish to decide if there is a feasible schedule.

**Theorem 1 (Rounding Theorem).** *Let  $P = (p_{ij}) \in \mathbb{Z}_+^{m \times n}$ ,  $\vec{d} = (d_1, \dots, d_m) \in \mathbb{Z}_+^m$  and  $t \in \mathbb{Z}_+$ . If the linear program  $\text{LP}(P, \vec{d}, t)$ , given by*

$$\begin{aligned} \sum_{i \in M_j(t)} x_{ij} &= 1 \quad \text{for } j = 1, \dots, n, \\ \sum_{j \in J_i(t)} p_{ij} x_{ij} &\leq d_i \quad \text{for } i = 1, \dots, m, \\ x_{ij} &\geq 0 \quad \text{for } j \in J_i(t), \quad i = 1, \dots, m, \end{aligned}$$

*has a feasible solution, then any vertex  $\tilde{x}$  of this polytope can be rounded to a feasible solution  $\bar{x}$  of the integer program  $\text{IP}(P, \vec{d}, t)$ , given by*

$$\begin{aligned} \sum_{i \in M_j(t)} x_{ij} &= 1 \quad \text{for } j = 1, \dots, n, \\ \sum_{j \in J_i(t)} p_{ij} x_{ij} &\leq d_i + t \quad \text{for } i = 1, \dots, m, \\ x_{ij} &\in \{0, 1\} \quad \text{for } j \in J_i(t), \quad i = 1, \dots, m, \end{aligned}$$

*and this rounding can be done in polynomial time.*

**Proof.** Let  $v$  denote the number of variables in the linear program  $\text{LP}(P, \vec{d}, t)$ . This polyhedron is defined by  $v + m + n$  constraints and is contained in the unit hypercube. Each vertex of such a *pointed* polyhedron is determined by  $v$  linearly independent rows of the constraint matrix such that each of these constraints is satisfied with equality (Schrijver, 1986). As a result, for any vertex  $\tilde{x}$  all but  $m + n$  of the variables must have value 0, and a straightforward counting argument shows that all non-integral variables correspond to at most  $m$  jobs. In the remainder of the proof, we first show a somewhat stronger structural property, and then use it to round the solution.

It will be convenient to associate the rows and columns of  $P$  with machines and jobs, as is true in our application. Suppose that  $\text{LP}(P, \vec{d}, t)$  is feasible and let  $\tilde{x}$  be a vertex of  $\text{LP}(P, \vec{d}, t)$ . For any feasible solution  $x$ , define a bipartite graph  $G(x) = (M, J, E)$ , where  $M = \{1, \dots, m\}$  and  $J = \{1, \dots, n\}$  correspond to the sets of machines and jobs, respectively, and  $E = \{(i, j) \mid x_{ij} > 0\}$ . We have already indicated that  $G(\tilde{x})$  has no more edges than nodes. We now show that each connected component of  $G(\tilde{x})$  has this property; that is,  $G(\tilde{x})$  is a *pseudoforest*. (This result was already stated in a slightly different form by Dantzig, 1963.)

Let  $C$  be a connected component of  $G(\tilde{x})$ . If  $M_C$  and  $J_C$  are the sets of machine and job nodes contained in  $C$ , let  $\tilde{x}_C$  denote the restriction of  $\tilde{x}$  to components  $\tilde{x}_{ij}$  such that  $i \in M_C$  and  $j \in J_C$ , and let  $\tilde{x}_{\bar{C}}$  denote the remaining components of  $\tilde{x}$ . For simplicity of notation, reorder the components so that  $\tilde{x} = (\tilde{x}_C, \tilde{x}_{\bar{C}})$ . Let  $P_C$  denote the restriction of  $P$  where the only remaining rows correspond to machines in  $M_C$  and the only remaining columns correspond to jobs in  $J_C$ . Similarly, define  $\vec{d}_C$  to be the restriction of  $\vec{d}$  to machines in  $M_C$ .

We first prove that  $\tilde{x}_C$  is an extreme point of  $\text{LP}(P_C, \vec{d}_C, t)$ . Suppose not; then there exist  $y_1$  and  $y_2$  such that  $\tilde{x}_C = \frac{1}{2}(y_1 + y_2)$ , where each  $y_i$  is a feasible solution of  $\text{LP}(P_C, \vec{d}_C, t)$ . But now,  $\tilde{x} = \frac{1}{2}((y_1, \tilde{x}_{\bar{C}}) + (y_2, \tilde{x}_{\bar{C}}))$  where each  $(y_i, \tilde{x}_{\bar{C}})$  is a feasible solution of  $\text{LP}(P, \vec{d}, t)$ , contradicting the choice of  $\tilde{x}$ .

Since  $\tilde{x}_C$  is an extreme point of  $\text{LP}(P, \vec{d}, t)$ , it follows that  $G(\tilde{x}_C)$  has no more edges than nodes. Since  $C = G(\tilde{x}_C)$  is connected,  $C$  is either a tree or a tree plus one additional edge, so that  $G(\tilde{x})$  is a pseudoforest.

We now use the fact that  $G(\tilde{x})$  is a pseudoforest to round the corresponding vertex  $\tilde{x}$ . Consider each edge  $(i, j)$  with  $\tilde{x}_{ij} = 1$ . For each such edge, we adopt this integral schedule for job  $j$  and set  $\tilde{x}_{ij} = 1$ . These jobs correspond to the job nodes of degree 1, so that by deleting all of these nodes we get a pseudoforest  $G'(\tilde{x})$  with the additional property that each job node has degree at least 2.

We show that  $G'(\tilde{x})$  has a matching that covers all of the job nodes. For each component that is a tree, root the tree at any node, and match each job node with any one of its children. (Note that each job node must have at least one child and that, since each node has at most one parent, no machine is matched with more than one job.) For each component that contains a cycle, take alternate edges of the cycle in the matching. (Note that the cycle must be of even length.) If the edges of the cycle are deleted, we get a collection of trees which we think of as rooted at

the node that had been contained in the cycle. For each job node that is not already matched, pair it with one of its children. This gives us the desired matching. If  $(i, j)$  is in the matching, set  $\bar{x}_{ij} = 1$ . Each remaining  $\bar{x}_{ij}$  that has not been assigned is set to 0.

It is straightforward to verify that  $\bar{x}$  is a feasible solution to  $\text{IP}(P, \vec{d}, t)$ . Each job has been scheduled on exactly one machine, so that

$$\sum_{i \in M_j(t)} \bar{x}_{ij} = 1$$

for  $j = 1, \dots, n$ . For each machine  $i$ , there is at most one job  $j$  such that  $\tilde{x}_{ij} < \bar{x}_{ij} = 1$ ; since  $p_{ij} \leq t$  for each  $x_{ij}$  in  $\text{IP}(P, \vec{d}, t)$ , we have

$$\sum_{j \in J_i(t)} p_{ij} \bar{x}_{ij} \leq \sum_{j \in J_i(t)} p_{ij} \tilde{x}_{ij} + t \leq d_i + t$$

for  $i = 1, \dots, m$ .  $\square$

### 3. Approximation algorithms

We are now ready to present approximation algorithms for the minimum makespan problem on unrelated parallel machines. A notion that will play a central role in these algorithms is that of a  $\rho$ -relaxed decision procedure. Consider the following decision version of the problem: given a matrix  $P$  of processing times and a deadline  $d$ , does there exist a schedule with makespan at most  $d$ ? On input  $(P, d)$ , an ordinary decision procedure would output ‘yes’ or ‘no’, depending on whether there was in fact such a schedule. A  $\rho$ -relaxed decision procedure outputs ‘no’ or ‘almost’; more precisely, on input  $(P, d)$ ,

- (i) it either outputs ‘no’ or produces a schedule with makespan at most  $\rho d$ , and
- (ii) if the output is ‘no’, then there is no schedule with makespan at most  $d$ .

Variations of the following lemma have been used in several recent results on approximation algorithms for scheduling problems (Hochbaum and Shmoys, 1987).

**Lemma 1.** *If there is a polynomial  $\rho$ -relaxed decision procedure for the minimum makespan problem on unrelated parallel machines, then there is a polynomial  $\rho$ -approximation algorithm for this problem.*

**Proof.** On input  $P$ , construct a greedy schedule, where each job is assigned to the machine on which it runs fastest. If the makespan of this schedule is  $t$ , then  $t$  is an upper bound on the optimum makespan, whereas  $t/m$  is a lower bound. Using these initial bounds, we run a binary search procedure. If  $u$  and  $l$  are the current upper and lower bounds, set  $d = \lfloor \frac{1}{2}(u + l) \rfloor$  and apply the  $\rho$ -relaxed decision procedure to  $(P, d)$ . If the answer is ‘almost’, then reset  $u$  to  $d$ , and otherwise reset  $l$  to  $d + 1$ , while storing the best solution obtained so far. When the upper and lower bounds are equal, output the best solution found.

It is easy to see that this procedure has the appropriate performance guarantee. A trivial inductive argument shows that  $l$  is always a lower bound on the optimum makespan and that the best solution encountered has makespan at most  $\rho u$ . Since  $u = l$  at termination, we get the desired bound. Furthermore, the algorithm clearly runs in polynomial time. The difference between the upper and lower bounds after  $k + \log m$  iterations is bounded by  $2^{-k}$  times the optimum. Thus after a polynomial number of iterations, the difference is less than 1 and the algorithm terminates.  $\square$

To obtain a polynomial 2-approximation algorithm is quite simple, given the rounding theorem and Lemma 1. We construct a 2-relaxed decision procedure for the decision version of the problem. Let  $(P, d)$  be a problem instance. Consider the linear program  $LP(P, \vec{d}, t)$  of the rounding theorem, with  $d_1 = \dots = d_m = t = d$ . If the instance is a 'yes' instance, then the schedule that completes by time  $d$  gives a feasible solution to the linear program: simply set  $x_{ij}$  to 1 if job  $j$  is assigned to machine  $i$  and 0 otherwise. In this case, the feasible region of the linear program is nonempty and pointed, and so it is possible to find a vertex  $\tilde{x}$  in polynomial time (Khachian, 1979; Grötschel, Lovász and Schrijver, 1988). Thus, if no vertex of  $LP(P, \vec{d}, t)$  is found, the instance must be a 'no' instance; otherwise, the procedure given in the proof of the rounding theorem produces a solution to the integer program  $IP(P, \vec{d}, t)$ . This 0-1 solution can be interpreted as a schedule in the obvious way, and it has makespan at most  $2d$ . Hence, the procedure is a 2-relaxed decision procedure. We have proved the following result.

**Theorem 2.** *There is a 2-approximation algorithm for the minimum makespan problem on unrelated parallel machines that runs in time bounded by a polynomial in the input size.*  $\square$

The analysis of the algorithm cannot be improved to yield a better worst-case performance bound. Consider the following instance. There are  $m^2 - m + 1$  jobs and  $m$  identical machines. The first job takes  $m$  time units on all machines, and all other jobs take one time unit on all machines. Clearly, the optimal schedule has makespan  $m$ : assign the first job to one machine and  $m$  of the remaining jobs to each of the other machines. No deadline less than  $m$  has a feasible fractional schedule. For  $d = m$ , suppose that the vertex of  $LP(P, \vec{d}, t)$  that is found corresponds to the schedule where one unit of the job of length  $m$  and  $m - 1$  unit length jobs are assigned to each machine. Rounding this fractional solution produces a schedule of length  $2m - 1$ . In this example, the choice of vertex has forced the inferior solution; no extension of the integral part of this solution will yield a better solution.

Even when  $m = 2$ , the algorithm can deliver a schedule of length arbitrarily close to twice the optimum. Consider the following instance. There are three jobs and two machines, with  $p_{11} = 1$ ,  $p_{21} = q$ ,  $p_{12} = q$ ,  $p_{22} = q + 1$ ,  $p_{13} = (3q + 2)/(2q + 1)$  and  $p_{23} = q + 1$ . The optimal schedule has makespan  $q + 1$ : jobs 2 and 3 use different machines, and job 1 uses machine 1. For any deadline  $d < q + 1$ , jobs 2 and 3 cannot

be scheduled on machine 2, and  $p_{12} + p_{13}$  exceeds  $q + 1$ , so  $LP(P, \vec{d}, t)$  will have no feasible solution. For  $d = q + 1$ , the vertex  $\tilde{x}$  returned could be given by  $\tilde{x}_{11} = 0$ ,  $\tilde{x}_{21} = 1$ ,  $\tilde{x}_{12} = \tilde{x}_{13} = (2q + 1)/(2q + 2)$  and  $\tilde{x}_{22} = \tilde{x}_{23} = 1/(2q + 2)$ . Job 1 has been integrally scheduled on machine 2, and jobs 2 and 3 must then be matched with machines 1 and 2. Either matching is possible, but in each case, the job assigned to machine 2 takes  $q + 1$  time units. This yields a schedule of length  $2q + 1$ . If we would not be restricted to a matching for the fractional jobs and assign both jobs 2 and 3 to machine 1, we obtain a schedule of length  $q + (3q + 2)/(2q + 1)$ , which is close to optimal.

As a second application of the rounding theorem, we give a polynomial approximation scheme for any fixed number  $m$  of machines. The running time of the procedure  $A_\varepsilon$ , which produces a schedule guaranteed to be within a factor  $1 + \varepsilon$  of the optimum, will be bounded by a function that is the product of  $(n + 1)^{m/\varepsilon}$  and a polynomial in the size of the input  $P$ . Given the fully polynomial approximation scheme of Horowitz and Sahni (1976), one may question the novelty of such a scheme. The significance of the new result lies in the fact that the space required by the old scheme is  $(nm/\varepsilon)^m$  whereas the new scheme uses space that is polynomial in both  $\log(1/\varepsilon)$  and  $m$  (and the input size).

Again, all we have to do is to construct a  $(1 + \varepsilon)$ -relaxed decision procedure for the decision version of the problem. Given  $(P, d)$ , the algorithm attempts to find solutions to  $(n + 1)^{m/\varepsilon}$  linear programs. If there is a schedule with makespan at most  $d$ , then one of these linear programs has a feasible solution, and this will correspond to a schedule with makespan at most  $(1 + \varepsilon)d$ . This suffices to guarantee the properties of a  $(1 + \varepsilon)$ -relaxed decision procedure.

For any schedule for the instance  $(P, d)$ , we classify the assignment of a job to a machine as either *long* or *short*, depending on whether or not the processing time in question is greater than  $\varepsilon d$ . No machine can handle  $1/\varepsilon$  or more long assignments before time  $d$ . Thus, for any instance there are less than  $(n + 1)^{m/\varepsilon}$  schedules of long assignments.

If the instance  $(P, d)$  has a feasible schedule, then this includes a partial schedule of long assignments (which may be empty). Suppose that for machine  $i$  the long assignments amount to a total processing time  $t_i$ , and thus the remaining jobs are completed within time  $d_i = d - t_i$ . If we then set  $t = \varepsilon d$ , we see that the linear program  $LP(P, \vec{d}, t)$  must once again have a feasible solution, so that we can apply the rounding theorem. The resulting integral solution yields a schedule of short assignments such that the total time taken by short assignments to machine  $i$  is at most  $d - t_i + \varepsilon d$ . Combining this with the schedule of long assignments, we get a schedule where the total time used by machine  $i$  is at most  $t_i + d - t_i + \varepsilon d = (1 + \varepsilon)d$ .

We try all possible schedules of long assignments in this way, computing the remaining available time on each machine and applying the rounding procedure. Either we conclude that the instance is a 'no' instance, or we produce a schedule with makespan at most  $(1 + \varepsilon)d$ . (Notice that if  $\varepsilon = 1$ , there are no long assignments and the algorithm reduces to the previous one.) We have shown the following result.



**Theorem 3.** *Let  $m$  be a fixed integer. There is a family  $\{A_\varepsilon\}$  of algorithms such that, for each  $\varepsilon > 0$ ,  $A_\varepsilon$  is a  $(1 + \varepsilon)$ -approximation algorithm for the minimum makespan problem on  $m$  unrelated parallel machines that requires time bounded by a polynomial in the input size and space bounded by a polynomial in  $m$ ,  $\log(1/\varepsilon)$ , and the input size.  $\square$*

An interesting open question is whether this result can be strengthened to give a *fully* polynomial approximation scheme for fixed values of  $m$ , where the space required is bounded by a polynomial in the input size,  $m$ , and  $1/\varepsilon$  (or, even better,  $\log(1/\varepsilon)$ ).

#### 4. Limits to approximation

We now present results which show that certain polynomial approximation algorithms cannot exist unless  $P = NP$ . To this end, we investigate the computational complexity of decision versions of our scheduling problem with small integral deadlines.

**Theorem 4.** *For the minimum makespan problem on unrelated parallel machines, the question of deciding if there exists a schedule with makespan at most 3 is NP-complete.*

We prove this result by a reduction from the 3-dimensional matching problem, which is known to be NP-complete:

#### 3-Dimensional matching.

*Instance:* Disjoint sets  $A = \{a_1, \dots, a_n\}$ ,  $B = \{b_1, \dots, b_n\}$ ,  $C = \{c_1, \dots, c_n\}$ , and a family  $F = \{T_1, \dots, T_m\}$  of triples with  $|T_i \cap A| = |T_i \cap B| = |T_i \cap C| = 1$  for  $i = 1, \dots, m$ .

*Question:* Does  $F$  contain a matching, i.e., a subfamily  $F'$  for which  $|F'| = n$  and  $\bigcup_{T_i \in F'} T_i = A \cup B \cup C$ ?

**Proof of Theorem 4.** Given an instance of the above 3-dimensional matching problem, we construct an instance of the scheduling problem with  $m$  machines and  $2n + m$  jobs. Machine  $i$  corresponds to the triple  $T_i$  for  $i = 1, \dots, m$ . There are  $3n$  ‘element jobs’ that correspond to the  $3n$  elements of  $A \cup B \cup C$  in the natural way. In addition, there are  $m - n$  ‘dummy jobs’. (If  $m < n$ , we construct some trivial ‘no’ instance of the scheduling problem.) Machine  $i$  corresponding to  $T_i = (a_j, b_k, c_l)$  can process each of the jobs corresponding to  $a_j$ ,  $b_k$  and  $c_l$  in one time unit and each other job in three time units. Note that the dummy jobs require three time units on each machine.

It is quite simple to show that there is a schedule with makespan at most 3 if and only if there is a 3-dimensional matching. Suppose there is a matching. For each  $T_i = (a_j, b_k, c_l)$  in the matching, schedule the element jobs corresponding to  $a_j$ ,  $b_k$  and  $c_l$  on machine  $i$ . Schedule the dummy jobs on the  $m - n$  machines corresponding to the triples that are not in the matching. This gives a schedule with makespan 3. Conversely, suppose that there is such a schedule. Each of the dummy jobs requires three time units on any machine and is thus scheduled by itself on some machine. Consider the set of  $n$  machines that are not processing dummy jobs. Since these are processing all of the  $3n$  element jobs, each of these jobs is processed in one time unit. Each three jobs that are assigned to one machine must therefore correspond to elements that form the triple corresponding to that machine. Since each element job is scheduled exactly once, the  $n$  triples corresponding to the machines that are not processing dummy jobs form a matching.  $\square$

As an immediate corollary of Theorem 4, we get the following result.

**Corollary 1.** *For every  $\rho < \frac{4}{3}$ , there does not exist a polynomial  $\rho$ -approximation algorithm for the minimum makespan problem on unrelated parallel machines, unless  $P = NP$ .  $\square$*

The technique employed in Theorem 4 can be refined to yield a stronger result.

**Theorem 5.** *For the minimum makespan problem on unrelated parallel machines, the question of deciding if there exists a schedule with makespan at most 2 is NP-complete.*

**Proof.** We again start from the 3-dimensional matching problem. We call the triples that contain  $a_j$  *triples of type  $j$* . Let  $t_j$  be the number of triples of type  $j$  for  $j = 1, \dots, n$ . As before, machine  $i$  corresponds to the triple  $T_i$  for  $i = 1, \dots, m$ . There are now only  $2n$  element jobs, corresponding to the  $2n$  elements of  $B \cup C$ . We refine the construction of the dummy jobs: there are  $t_j - 1$  dummy jobs of type  $j$  for  $j = 1, \dots, n$ . (Note that the total number of dummy jobs is  $m - n$ , as before.) Machine  $i$  corresponding to a triple of type  $j$ , say  $T_i = (a_j, b_k, c_l)$ , can process each of the element jobs corresponding to  $b_k$  and  $c_l$  in one time unit and each of the dummy jobs of type  $j$  in two time units; all other jobs require three time units on machine  $i$ .

Suppose there is a matching. For each  $T_i = (a_j, b_k, c_l)$  in the matching, schedule the element jobs corresponding to  $b_k$  and  $c_l$  on machine  $i$ . For each  $j$ , this leaves  $t_j - 1$  idle machines corresponding to triples of type  $j$  that are not in the matching; schedule the  $t_j - 1$  dummy jobs of type  $j$  on these machines. This completes a schedule with makespan 2. Conversely, suppose that there is such a schedule. Each dummy job of type  $j$  is scheduled on a machine corresponding to a triple of type  $j$ . Therefore, there is exactly one machine corresponding to a triple of type  $j$  that is not processing dummy jobs for  $j = 1, \dots, n$ . Each such machine is processing two element jobs in one time unit each. If the machine corresponds to a triple of type

$j$  and its two unit-time jobs correspond to  $b_k$  and  $c_l$ , then  $(a_j, b_k, c_l)$  must be the triple corresponding to that machine. Since each element job is scheduled exactly once, the  $n$  triples corresponding to the machines that are not processing dummy jobs form a matching.  $\square$

**Corollary 2.** *For every  $\rho < \frac{3}{2}$ , there does not exist a polynomial  $\rho$ -approximation algorithm for the minimum makespan problem on unrelated parallel machines, unless  $P = NP$ .  $\square$*

## 5. Restricted processing times

We conclude this paper with a few remarks about special cases of our scheduling problem in which the number of different processing times is bounded. If all  $p_{ij} = 1$ , the problem is clearly solvable in polynomial time, and even if all  $p_{ij} \in \{1, \infty\}$ , the problem can be solved by bipartite cardinality matching. Theorem 6 shows that, if all  $p_{ij} \in \{1, 2\}$ , the problem is still solvable in polynomial time by matching techniques.

**Theorem 6.** *The minimum makespan problem on unrelated parallel machines is solvable in polynomial time in the case that all  $p_{ij} \in \{1, 2\}$ .*

**Proof.** The problem of deciding if there exists a schedule with makespan at most  $d$  can be transformed into the following problem, which is known to be solvable in polynomial time by matching techniques (Schrijver, 1983): Given a bipartite graph  $G = (S, T, E)$ , find a subgraph with a maximum number of edges in which each node in  $S$  has degree 0 or 2 and each node in  $T$  has degree 1.

In the case  $d = 2k$ , construct  $G$  as follows: there are  $k$  nodes in  $S$  for each machine; there is one node in  $T$  for each job; and edges correspond to unit processing times. Now solve the above problem on  $G$ . Construct a partial schedule of unit-time assignments corresponding to the edges in the subgraph, and try to extend it to a complete schedule by assigning the remaining jobs for two time units each while respecting the deadline. If a schedule with makespan  $2k$  exists, this procedure will find such a schedule.

In the case  $d = 2k - 1$ , create an additional unit-time job for each machine, which requires two time units on any other machine, and apply the above procedure for  $d = 2k$ . It follows from a straightforward interchange argument that the original problem has a schedule with makespan  $2k - 1$  if and only if the modified problem has a schedule with makespan  $2k$ .  $\square$

Our final theorem implies that all other cases with a fixed number of processing times cannot be solved in polynomial time, unless  $P = NP$ .

**Theorem 7.** *The minimum makespan problem on unrelated parallel machines is NP-hard in the case that all  $p_{ij} \in \{p, q\}$  with  $p < q$ ,  $2p \neq q$ .*

**Proof.** We assume without loss of generality that  $p$  and  $q$  are relatively prime. Recall that the result has already been proved in Theorem 4 for the case that all  $p_{ij} \in \{1, 3\}$  by a reduction from 3-dimensional matching. We now reduce from  $q$ -dimensional matching. Given a matching instance with  $m$   $q$ -tuples over a ground set of  $qn$  elements, we construct a scheduling instance with  $qn$  element jobs,  $p(m - n)$  dummy jobs and  $m$  machines. An element job can be processed in  $p$  time units by each machine that corresponds to a tuple containing the corresponding element; all other processing times are  $q$  time units. It is trivial to see that for each matching there is a schedule with makespan  $pq$ . To prove the opposite implication, no schedule with makespan  $pq$  can have idle time, and it now follows from an easy number theoretic argument that each machine processes either  $q$  element jobs of length  $p$  or  $p$  dummy jobs of length  $q$ .  $\square$

## Acknowledgements

The research of the first author was supported in part by the National Science Foundation under grant MCS-83-11422. The research by the second author was supported in part by the National Science Foundation under grant ECS-85-01988 and a Presidential Young Investigator award with matching support from IBM and Sun, and by Air Force contract AFOSR-86-0078. The research by the third author was supported in part by the National Science Foundation under grant MCS-81-20790 and by the Hungarian National Foundation for Scientific Research under grant 1812.

## References

- R. Aharoni, P. Erdős and N. Linial, "Dual integer linear programs and the relationship between their optima," *Proceedings of the 17th Annual ACM Symposium on Theory of Computing* (1985) 476–483.
- J.J. Bartholdi III, "A guaranteed-accuracy round-off algorithm for cyclic scheduling and set covering," *Operations Research* 29 (1981) 501–510.
- J.J. Bartholdi III, J.B. Orlin and H.D. Ratliff, "Cyclic scheduling via integer programs with circular ones," *Operations Research* 28 (1980) 1074–1085.
- S. Baum and L.E. Trotter Jr., "Integer rounding for polymatroid and branching optimization problems," *SIAM Journal on Algebraic and Discrete Methods* 2 (1981) 416–425.
- V. Chvátal, "A greedy heuristic for the set-covering problem," *Mathematics of Operations Research* 4 (1979) 233–235.
- G.B. Dantzig, *Linear Programming and Extensions* (Princeton University Press, Princeton, NJ, 1963).
- E. Davis and J.M. Jaffe, "Algorithms for scheduling tasks on unrelated processors," *Journal of the Association for Computing Machinery* 28 (1981) 721–736.
- M.R. Garey and D.S. Johnson, "Complexity results for multiprocessor scheduling under resource constraints," *SIAM Journal on Computing* 4 (1975) 397–411.

- M.R. Garey and D.S. Johnson, "Strong NP-completeness results: motivation, examples and implications," *Journal of the Association for Computing Machinery* 25 (1978) 499–508.
- M.R. Garey and D.S. Johnson, *Computers and Intractability: a Guide to the Theory of NP-Completeness* (Freeman, San Francisco, CA, 1979).
- T. Gonzalez, O.H. Ibarra and S. Sahni, "Bounds for LPT schedules on uniform processors," *SIAM Journal on Computing* 6 (1977) 155–166.
- R.L. Graham, "Bounds for certain multiprocessing anomalies," *Bell System Technological Journal* 45 (1966) 1563–1581.
- R.L. Graham, "Bounds on multiprocessing timing anomalies," *SIAM Journal on Applied Mathematics* 17 (1969) 416–429.
- R.L. Graham, E.L. Lawler, J.K. Lenstra and A.H.G. Rinnooy Kan, "Optimization and approximation in deterministic sequencing and scheduling: a survey," *Annals of Discrete Mathematics* 5 (1979) 287–326.
- M. Grötschel, L. Lovász and A. Schrijver, *Geometric Algorithms and Combinatorial Optimization* (Springer, Berlin, 1988).
- D.S. Hochbaum and D.B. Shmoys, "Using dual approximation algorithms for scheduling problems: practical and theoretical results," *Journal of the Association for Computing Machinery* 34 (1987) 144–162.
- D.S. Hochbaum and D.B. Shmoys, "A polynomial approximation scheme for machine scheduling on uniform processors: using the dual approximation approach," *SIAM Journal on Computing* 17 (1988) 539–551.
- E. Horowitz and S. Sahni, "Exact and approximate algorithms for scheduling nonidentical processors," *Journal of the Association for Computing Machinery* 23 (1976) 317–327.
- L.G. Khachian, "A polynomial time algorithm in linear programming," *Soviet Mathematics Doklady* 20 (1979) 191–194.
- L. Lovász, "On the ratio of optimal and fractional covers," *Discrete Applied Mathematics* 13 (1983) 383–390.
- O.M.-C. Marcotte, *Topics in Combinatorial Packing and Covering*, Ph.D. Thesis, School of Operations Research and Industrial Engineering, Cornell University (Ithaca, NY, 1983).
- K. Numata, "Approximate and exact algorithms for scheduling independent tasks on unrelated processors," *Journal of the Operations Research Society of Japan* 31 (1988) 61–81.
- C.N. Potts, "Analysis of a linear programming heuristic for scheduling unrelated parallel machines," *Discrete Applied Mathematics* 10 (1985) 155–164.
- P. Raghavan, "Probabilistic construction of deterministic algorithms: approximating packing integer programs," *Proceedings of the 27th Annual IEEE Symposium on Foundations of Computer Science* (1986) 10–18.
- P. Raghavan and C.D. Thompson, "Provably good routing in graphs: regular arrays," *Proceedings of the 17th Annual ACM Symposium on Theory of Computing* (1985) 79–87.
- S. Sahni, "Algorithms for scheduling independent tasks," *Journal of the Association for Computing Machinery* 23 (1976) 116–127.
- A. Schrijver, "Min–Max results in combinatorial optimization," in: A. Bachem, M. Grötschel and B. Korte, eds., *Mathematical Programming: The State of the Art—Bonn 1982* (Springer, Berlin, 1983) pp. 439–500.
- A. Schrijver, *Theory of Linear and Integer Programming* (Wiley, Chichester, 1986).