

# Planning and Scheduling to Minimize Tardiness

J. N. Hooker

Carnegie Mellon University  
john@hooker.tepper.cmu.edu

**Abstract.** We combine mixed integer linear programming (MILP) and constraint programming (CP) to minimize tardiness in planning and scheduling. Tasks are allocated to facilities using MILP and scheduled using CP, and the two are linked via logic-based Benders decomposition. We consider two objectives: minimizing the number of late tasks, and minimizing total tardiness. Our main theoretical contribution is a relaxation of the cumulative scheduling subproblem, which is critical to performance. We obtain substantial computational speedups relative to the state of the art in both MILP and CP. We also obtain much better solutions for problems that cannot be solved to optimality.

We address a planning and scheduling problem that occurs frequently in manufacturing and supply chain contexts. Tasks must be assigned to facilities and scheduled on each facility subject to release dates and due dates. Tasks assigned to a given facility may run in parallel if desired, subject to a resource constraint (cumulative scheduling). We consider two objectives: minimizing the number of late tasks, and minimizing total tardiness.

The problem can be formulated entirely as a constraint programming (CP) problem or a mixed integer/linear programming (MILP) problem. However, these models are hard to solve. By linking CP and MILP in a hybrid method, we obtain significant speedups relative to the state of the art in both MILP and CP. The linkage is achieved by logic-based Benders decomposition. The facility assignment problem becomes the master problem and is solved by MILP, while the scheduling problem becomes the subproblem (slave problem) and is solved by CP.

The primary theoretical contribution of this paper is a linear relaxation of the cumulative scheduling subproblem. We find that including such a relaxation in the master problem is essential to the success of the Benders method.

We solve problem instances in which tasks have the same release date and different due dates, although the method is valid for different release dates as well. We obtain substantial speedups on nearly all instances relative to MILP (as represented by CPLEX), which in turn is generally faster than CP (as represented by the ILOG Scheduler). On larger instances, the hybrid method generally achieves speedups of two or three orders of magnitude when minimizing the number of late tasks, and it solves significantly more problems to optimality.

There is a lesser but still significant speedup when minimizing total tardiness, and even when the hybrid method cannot obtain provably optimal solutions, it obtains much better solutions than provided by MILP in the same time period.

## 1 Previous Work

Logic-based Benders decomposition, which generalizes the classical Benders decomposition method [1, 3], was introduced in [8] for purposes of logic circuit verification. The idea was later formally developed in [4] and applied to 0-1 programming in [7].

The application of logic-based Benders to planning and scheduling was proposed in [4]. Jain and Grossmann [10] successfully applied such a method to minimum-cost planning and scheduling problems in which the subproblems are disjunctive scheduling problems, where tasks must run one at a time, rather than cumulative scheduling problems. The Benders cuts are particularly simple in this case because the subproblem is a feasibility problem rather than an optimization problem. Thorsteinsson [11] improved on these results using a “branch-and-check” method suggested in [4].

It is less obvious how to define Benders cuts when the subproblem is an optimization problem. We showed in [5, 6] how to derive effective Benders cuts for at least one such case, minimum makespan problems. The cuts are valid for cumulative as well as disjunctive scheduling, provided all tasks have the same release date. Computational tests showed the hybrid method to be 100 to 1000 times faster than MILP or CP when all tasks have the same deadline.

In this paper we address minimum tardiness problems, in which the subproblem is again an optimization problem. We obtain effective cuts by repeatedly solving the subproblem with slightly different task assignments. The idea is related to finding “minimal conflict sets” of tasks, or small sets of tasks that create infeasibility when assigned to a particular facility. Cambazard et al. [2] applied such an approach to real-time scheduling of computing resources. Here we develop cuts for an optimization rather than a feasibility subproblem.

As observed in [5, 11], the success of hybrid methods in planning and scheduling relies on including a relaxation of the scheduling subproblem in the master problem. We find that deriving a useful relaxation requires deeper analysis when minimizing total tardiness than when minimizing cost or makespan. A relaxation of the cumulative scheduling problem is presented in [9], but it is expressed in terms of the start time variables, rather than the assignment variables as required for the Benders master problem. We derive here a very different relaxation in terms of 0-1 assignment variables, which is suitable for the MILP master problem.

## 2 The Problem

The planning and scheduling problem may be defined as follows. Each task  $j \in \{1, \dots, n\}$  is to be assigned to a facility  $i \in \{1, \dots, m\}$ , where it consumes

processing time  $p_{ij}$  and resources at the rate  $c_{ij}$ . Each task  $j$  has release time  $r_j$  and due date  $d_j$ . The tasks assigned to facility  $i$  must be given start times  $s_j$  in such a way that the total rate of resource consumption on facility  $i$  is never more than  $C_i$  at any given time. If  $x_j$  is the facility assigned to task  $j$ , the problem may be written

$$\begin{aligned} & \text{minimize } g(x, s) \\ & \text{subject to } r_j \leq s_j, \quad \text{all } j \quad (a) \\ & \quad \sum_{j \in J_{it}(x)} c_{ij} \leq C_i, \quad \text{all } i, t \quad (b) \end{aligned} \tag{1}$$

where  $x_j, s_j$  are the variables and  $J_{it}(x) = \{j \mid x_j = i, s_j \leq t \leq s_j + p_{ij}\}$  is the set of tasks underway at time  $t$  in facility  $i$ .

Precedence constraints may be imposed on tasks that are assigned to the same machine. Thus one may require that tasks  $j$  and  $k$  be scheduled on the same facility, and that task  $j$  precede  $k$ , by writing the constraints  $x_j = x_k$  and  $s_j + p_{x_j j} \leq s_k$ .

We investigate two objective functions:

- *number of late tasks*, given by  $g(x, s) = \sum_j \delta(s_j + p_{x_j j} - d_j)$ , where  $\delta(\alpha)$  is 1 if  $\alpha > 0$  and 0 otherwise.
- *total tardiness*, given by  $g(x, s) = \sum_j (s_j + p_{x_j j} - d_j)^+$ , where  $\alpha^+$  is  $\alpha$  if  $\alpha > 0$  and 0 otherwise.

### 3 Constraint Programming Formulation

A CP formulation of the problem can be written

$$\begin{aligned} & \text{minimize } g(x, s) \\ & \text{subject to } r_j \leq s_j, \quad \text{all } j \quad (2) \\ & \quad \text{cumulative}((s_j | x_j = i), (p_{ij} | x_j = i), (c_{ij} | x_j = i), C_i), \quad \text{all } i \end{aligned}$$

where  $(s_j | x_j = i)$  denotes the tuple of start times for tasks assigned to facility  $i$ . When minimizing the number of late tasks,  $g(x, s) = \sum_j L_j$  where  $L_j$  is binary, and the constraint  $(s_j + p_{x_j j} > d_j) \Rightarrow (L_j = 1)$  is added for each  $j$ . When minimizing total tardiness,  $g(x, s) = \sum_j T_j$ , and the constraints  $T_j \geq s_j + p_{x_j j} - d_j$  and  $T_j \geq 0$  are added for each  $j$ .

For purposes of computational testing we formulated (2) using the modeling language of OPL Studio. We used the `assignAlternatives` and `setTimes` search options specify a branching method that results in substantially better performance than the default method.

### 4 Mixed Integer Programming Formulation

The most straightforward MILP formulation discretizes time and enforces the resource capacity constraint at each discrete time. Let the 0-1 variable  $x_{ijt} = 1$

if task  $j$  starts at discrete time  $t$  on facility  $i$ . The formulation for minimizing the number of late tasks is

$$\begin{aligned}
\min \quad & \sum_j L_j \\
\text{subject to} \quad & NL_j \geq \sum_i (t + p_{ij})x_{ijt} - d_j, \quad \text{all } j, t \quad (a) \\
& \sum_{it} x_{ijt} = 1, \quad \text{all } j \quad (b) \quad (3) \\
& \sum_j \sum_{t' \in T_{ijt}} c_{ij}x_{ijt'} \leq C_i, \quad \text{all } i, t \quad (c) \\
& x_{ijt} = 0, \quad \text{all } j, t \text{ with } t < r_j \text{ or } t > N - p_{ij} \quad (d)
\end{aligned}$$

where each  $x_{ijt}$  and each  $L_j$  is a 0-1 variable. Also  $N$  is the number of discrete times (starting with  $t = 0$ ), and  $T_{ijt} = \{t' \mid t - p_{ij} < t' \leq t\}$  is the set of discrete times at which a task  $j$  in progress on facility  $i$  at time  $t$  might start processing. Constraint (b) ensures that each task starts once on one facility, (c) enforces the resource limit, and (d) the time windows. The minimum tardiness problem replaces the objective function with  $\sum_j T_j$  and constraint (a) with

$$T_j \geq \sum_i (t + p_{ij})x_{ijt} - d_j, \quad T_j \geq 0, \quad \text{all } j, t$$

We also investigated a smaller discrete event model suggested by [12], which uses continuous time. However, it proved much harder to solve than (3).

## 5 Hybrid Method for Minimizing Late Tasks

The Benders approach formulates a *master problem* that assigns tasks to facilities and a *subproblem* that schedules the tasks assigned to each facility. We write the master problem using an MILP model that minimizes the number of late tasks. In iteration  $h$  of the Benders algorithm, the master problem is

$$\begin{aligned}
\text{minimize} \quad & L \\
\text{subject to} \quad & \sum_i x_{ij} = 1, \quad \text{all } j \quad (a) \\
& \text{Benders cuts generated in iterations } 1, \dots, h-1 \quad (b) \\
& \text{relaxation of subproblem} \quad (c)
\end{aligned} \quad (4)$$

Here the binary variable  $x_{ij}$  is 1 when task  $j$  is assigned to facility  $i$ . The Benders cuts and relaxation will be described shortly.

Once an assignment  $\bar{x}_{ij}$  of tasks to facilities is determined by solving the master problem, a cumulative scheduling subproblem is solved by CP. The sub-

problem decouples into a separate scheduling problem on each facility  $i$ :

$$\begin{aligned}
& \text{minimize} && \sum_{j \in J_{hi}} L_j \\
& \text{subject to} && (s_j + p_{ij} > d_j) \Rightarrow (L_j = 1), \quad \text{all } j \in J_{hi} \\
& && r_j \leq s_j, \quad \text{all } j \in J_{hi} \\
& && \text{cumulative}((s_j|j \in J_{hi}), (p_{ij}|j \in J_{hi}), (c_{ij}|j \in J_{hi}))
\end{aligned} \tag{5}$$

where  $J_{hi}$  is the set of tasks for which  $\bar{x}_{ij} = 1$  (i.e., the tasks assigned to facility  $i$  in the master problem solution). If  $L_{hi}^*$  is the optimal value of (5), then  $\sum_i L_{hi}^*$  is the minimum number of late tasks across all facilities.

At this point we know that whenever the tasks in  $J_{hi}$  (perhaps among others) are assigned to facility  $i$ , the number of late tasks on facility  $i$  is at least  $L_{hi}^*$ . This allows us to write a valid lower bound  $\underline{L}_{hi}$  on the number of late tasks in facility  $i$  for any assignment of tasks to machines. Since  $x_{ij} = 0$  when task  $j$  is not assigned to facility  $i$ , we have

$$\begin{aligned}
\underline{L}_{hi} &\geq L_{hi}^* - L_{hi}^* \sum_{j \in J_{hi}} (1 - x_{ij}), \quad \text{all } i \\
\underline{L}_{hi} &\geq 0, \quad \text{all } i
\end{aligned} \tag{6}$$

By summing over all facilities, we have a lower bound on the total number  $L$  of late tasks:

$$L \geq \sum_i \underline{L}_{hi} \tag{7}$$

The inequality (7), together with (6), provides a *Benders cut* for iteration  $h$ . The cut says that the number of late tasks will be at least the number obtained in the subproblem unless a different assignment of tasks to facilities is used.

In iteration  $h$ , the Benders cuts (b) in the master problem (4) consist of inequalities (6)–(7) obtained in iterations  $1, \dots, h-1$ . The algorithm terminates when the optimal value of the master problem equals the optimal value of the subproblem in the previous iteration. At any point in the algorithm, a feasible solution of the subproblem is a feasible solution of the original problem, and the optimal value of the master problem is a lower bound on the optimal value of the original problem.

Unfortunately the Benders cuts (6)–(7) are weak and do not perform well in practice. The cuts can be strengthened by identifying, for each facility  $i$ , a smaller set  $J_{hi}$  of tasks that result in the same number of late tasks. One way to do this is to track which tasks actually play a role in the determining the minimum number of late tasks, as suggested in [6]. However, since this information is not available from commercial CP solvers, the information must be obtained indirectly by repeatedly solving subproblems with different assignments of tasks to facilities.

The following approach was found to yield effective cuts with a modest amount of computation. Let  $L_i(J)$  be the minimum number of late tasks on facility  $i$  when the tasks in  $J$  are assigned to facility  $i$ . First identify a set

Let  $J_{hi}^0 = J_{hi}$ .  
For all  $j \in J_{hi}$ : if  $L_i(J_{hi}^0 \setminus \{j\}) = L_{hi}^*$  then let  $J_{hi}^0 = J_{hi}^0 \setminus \{j\}$ .  
Let  $J_{hi}^1 = J_{hi}^0$ .  
For all  $j \in J_{hi}^0$ : if  $L_i(J_{hi}^1 \setminus \{j\}) = L_{hi}^* - 1$  then let  $J_{hi}^1 = J_{hi}^1 \setminus \{j\}$ .

**Fig. 1.** Algorithm for generating Benders cuts when minimizing the number of late tasks.

$J_{hi}^0 \subseteq J_{hi}$  of tasks that, when assigned to facility  $i$ , result in a minimum of  $L_{hi}^*$  late tasks; that is, a set  $J_{hi}^0$  such that  $L_i(J_{hi}^0) = L_{hi}^*$ . This is done via the simple greedy algorithm in Fig. 1. Then identify a set  $J_{hi}^1 \subseteq J_{hi}^0$  of tasks such that  $L_i(J_{hi}^1) = L_{hi}^* - 1$ , again using the algorithm of Fig. 1. The inequalities (6) can now be replaced by the generally stronger inequalities

$$\begin{aligned} \underline{L}_{hi} &\geq L_{hi}^* - L_{hi}^* \sum_{j \in J_{hi}^0} (1 - x_{ij}), \text{ all } i \\ \underline{L}_{hi} &\geq L_{hi}^* - 1 - L_{hi}^* \sum_{j \in J_{hi}^1} (1 - x_{ij}), \text{ all } i \\ \underline{L}_{hi} &\geq 0, \text{ all } i \end{aligned} \tag{8}$$

These cuts remain valid for any set of additional constraints that may be added to the subproblems.

It is straightforward to relax the subproblem when minimizing the number of late tasks. (It will be harder when minimizing total tardiness.) Let  $J(t_1, t_2)$  be the set of tasks whose time windows are contained in  $[t_1, t_2]$ . Thus  $J(t_1, t_2) = \{j \mid [r_j, d_j] \subseteq [t_1, t_2]\}$ . When executed on facility  $i$ , these tasks span a time interval of at least

$$M = \frac{1}{C_i} \sum_{j \in J(t_1, t_2)} c_{ij} p_{ij} \tag{9}$$

If  $M > t_2 - t_1$  then at least one task is late, and in fact the number of late tasks on facility  $i$  is at least

$$\frac{M - (t_2 - t_1)}{\max_{j \in J(t_1, t_2)} \{p_{ij}\}} \tag{10}$$

rounded up to the nearest integer.

Define  $\bar{r}_1, \dots, \bar{r}_{n_r}$  to be the distinct values among the release times  $r_1, \dots, r_n$  in increasing order, and similarly for  $\bar{d}_1, \dots, \bar{d}_{n_d}$ . Then from (9) and (10) we have the following relaxation:

$$\begin{aligned} L &\geq \sum_i L_i \\ L_i &\geq \frac{\frac{1}{C_i} \sum_{\ell \in J(\bar{r}_j, \bar{d}_k)} c_{i\ell} p_{i\ell} x_{i\ell} - (\bar{d}_k - \bar{r}_j)}{\max_{\ell \in J(\bar{r}_j, \bar{d}_j)} \{p_{i\ell}\}}, \quad j = 1, \dots, n_r, \quad k = 1, \dots, n_d, \text{ all } i \end{aligned}$$

which becomes (c) in the master problem (4).

## 6 Hybrid Method for Minimizing Total Tardiness

In iteration  $h$  of the Benders method, the master problem for minimizing total tardiness is

$$\begin{aligned}
& \text{minimize} && T \\
& \text{subject to} && \sum_i x_{ij} = 1, \text{ all } j && (a) \\
& && \text{Benders cuts for iterations } 1, \dots, h-1 && (b) \\
& && \text{relaxation of subproblem} && (c)
\end{aligned} \tag{11}$$

The subproblem again decouples into a cumulative scheduling problem for each facility  $i$ :

$$\begin{aligned}
& \text{minimize} && \sum_{j \in J_i} T_j \\
& \text{subject to} && T_j \geq s_j + p_{ij} - d_j, \text{ all } j \in J_i && (12) \\
& && r_j \leq s_j, \text{ all } j \in J_i \\
& && \text{cumulative}((s_j | j \in J_i), (p_{ij} | j \in J_i), (c_{ij} | j \in J_i))
\end{aligned}$$

We found the following scheme to generate effective Benders cuts. As before let  $J_{hi}$  be a set of tasks assigned to facility  $i$  in iteration  $h$ , and let  $T_{hi}^*$  be the resulting minimum tardiness on facility  $i$ . Let  $T_i(J)$  be the minimum tardiness on facility  $i$  that results when the tasks in  $J$  are assigned to facility  $i$ , so that  $T_i(J_{hi}) = T_{hi}^*$ . Let  $Z_{hi}$  be the set of tasks in  $J_{hi}$  that can be removed, one at a time, without reducing the minimum tardiness. That is,

$$Z_{hi} = \{j \in J_{hi} \mid T_{hi}(J_{hi} \setminus \{j\}) = T_{hi}^*\}$$

Finally, let  $T_{hi}^0$  be the minimum tardiness that results from removing the tasks in  $Z_{hi}$  all at once, so that  $T_{hi}^0 = T_i(J_{hi} \setminus Z_{hi})$ . Thus any or all tasks in  $Z_{hi}$  can be removed from facility  $i$  without reducing the minimum tardiness below  $T_{hi}^0$ . This yields the following Benders cuts in iteration  $h$ :

$$\begin{aligned}
T &\geq T_{hi}^0 - T_{hi}^0 \sum_{j \in J_{hi} \setminus Z_{hi}} (1 - x_{ij}), \text{ all } i \\
T &\geq T_{hi}^* - T_{hi}^* \sum_{j \in J_{hi}} (1 - x_{ij}), \text{ all } i
\end{aligned} \tag{13}$$

The second cut is redundant and can be eliminated for a given  $h, i$  when  $T_{hi}^0 = T_{hi}^*$ . This in fact substantially reduces the size of master problem, since computational testing suggests that  $T_{hi}^0 = T_{hi}^*$  very often.

These cuts are again valid for any set of additional constraints that may be added to the subproblem.

## 7 Relaxation for Minimizing Total Tardiness

Our relaxation of the minimum tardiness scheduling subproblem has two parts. The first and simpler part is similar to the relaxation obtained for minimizing the number of late tasks. It is based on the following lemma. Recall that  $J(t_1, t_2)$  is the set of jobs with time windows between  $t_1$  and  $t_2$ .

**Lemma 1.** *Consider a minimum total tardiness problem in which tasks  $j = 1, \dots, n$  with time windows  $[r_j, d_j]$  are scheduled on a single facility  $i$ , where  $\min_j \{r_j\} = 0$ . The total tardiness incurred by any feasible solution is bounded below by*

$$\left( \frac{1}{C_i} \sum_{j \in J(0, d_k)} p_{ij} c_{ij} - d_k \right)^+$$

for each  $k = 1, \dots, n$ .

*Proof.* For any  $k$ , the last scheduled task in the set  $J(0, d_k)$  can finish no earlier than time  $t = \frac{1}{C_i} \sum_{j \in J(0, d_k)} p_{ij} c_{ij}$ . Since the last task has due date no later than  $d_k$ , its tardiness is no less than  $(t - d_k)^+$ . Thus total tardiness is no less than  $(t - d_k)^+$ .

This gives rise to a relaxation consisting of

$$T \geq \frac{1}{C_i} \sum_{j \in J(0, d_k)} p_{ij} c_{ij} x_{ij} - d_k, \quad \text{all } i, k \quad (14)$$

and  $T \geq 0$ .

The second part of the relaxation can be developed on basis of the following lemma. For each facility  $i$  let  $\pi_i$  be a permutation of  $\{1, \dots, n\}$  such that  $p_{i\pi_i(1)} c_{i\pi_i(1)} \leq \dots \leq p_{i\pi_i(n)} c_{i\pi_i(n)}$ .

**Lemma 2.** *Consider a minimum tardiness problem in which tasks  $1, \dots, n$  with time windows  $[r_j, d_j]$  are scheduled on a single facility  $i$ . Assume  $\min_j \{r_j\} = 0$  and index the tasks so that  $d_1 \leq \dots \leq d_n$ . Then the total tardiness  $T$  of any feasible solution is bounded below by  $\underline{T} = \sum_{k=1}^n \underline{T}_k$ , where*

$$\underline{T}_k = \left( \frac{1}{C_i} \sum_{j=1}^k p_{i\pi_i(j)} c_{i\pi_i(j)} - d_k \right)^+, \quad k = 1, \dots, n$$

*Proof.* Consider any feasible solution of the one-facility minimum tardiness problem, in which tasks  $1, \dots, n$  are respectively scheduled at times  $t_1, \dots, t_n$ . Thus

$$T = \sum_{k=1}^n (t_k + p_{ik} - d_k)^+ \quad (15)$$



Let  $\sigma_0(1), \dots, \sigma_0(n)$  be the order in which tasks are scheduled in this solution, so that  $t_{\sigma_0(1)} \leq \dots \leq t_{\sigma_0(n)}$ . For an arbitrary permutation  $\sigma$  of  $\{1, \dots, n\}$  let

$$\underline{T}_k(\sigma) = \left( \frac{1}{C_i} \sum_{j=1}^k p_{i\pi_i(j)} c_{i\pi_i(j)} - d_{\sigma(k)} \right)^+ \quad (16)$$

and  $\underline{T}(\sigma) = \sum_{k=1}^n \underline{T}_k(\sigma)$ .

We show first that  $T \geq \underline{T}(\sigma_0)$ . Since  $\sigma_0$  is a permutation we can write (15) as

$$T = \sum_{k=1}^n (t_{\sigma_0(k)} + p_{i\sigma_0(k)} - d_{\sigma_0(k)})^+$$

We observe that

$$t_{\sigma_0(k)} + p_{i\sigma_0(k)} \geq \frac{1}{C_i} \sum_{j=1}^k p_{i\sigma_0(j)} c_{i\sigma_0(j)} \geq \frac{1}{C_i} \sum_{j=1}^k p_{i\pi_i(j)} c_{i\pi_i(j)}$$

where the first inequality is based on the areas required by tasks, and the second inequality is due to the definition of  $\pi_i$ . From this and (16) it follows that  $T \geq \underline{T}(\sigma_0)$ .

Now suppose a bubble sort is performed on the integers  $\sigma_0(1), \dots, \sigma_0(n)$  so as to put them in increasing order, and let  $\sigma_0, \dots, \sigma_P$  be the resulting series of permutations. Thus  $(\sigma_P(1), \dots, \sigma_P(n)) = (1, \dots, n)$ , and  $\sigma_{p+1}$  is obtained from  $\sigma_p$  by swapping two adjacent terms  $\sigma_p(k)$  and  $\sigma_p(k+1)$ , where  $\sigma_p(k) > \sigma_p(k+1)$ . This means  $\sigma_p$  and  $\sigma_{p+1}$  are the same except that  $\sigma_{p+1}(k) = \sigma_p(k+1)$  and  $\sigma_{p+1}(k+1) = \sigma_p(k)$ . Since  $T^* \geq \underline{T}(\sigma_0)$  and  $\underline{T}(\sigma_P) = \underline{T}$ , to prove the theorem it suffices to show  $\underline{T}(\sigma_0) \geq \dots \geq \underline{T}(\sigma_P)$ .

Thus we consider any two adjacent permutations  $\sigma_p, \sigma_{p+1}$  and show that  $\underline{T}(\sigma_p) \geq \underline{T}(\sigma_{p+1})$ . We observe that

$$\begin{aligned} \underline{T}(\sigma_p) &= \sum_{j=1}^{k-1} \underline{T}_j(\sigma_p) + \underline{T}_k(\sigma_p) + \underline{T}_{k+1}(\sigma_p) + \sum_{j=k+2}^n \underline{T}_j(\sigma_p) \\ \underline{T}(\sigma_{p+1}) &= \sum_{j=1}^{k-1} \underline{T}_j(\sigma_p) + \underline{T}_k(\sigma_{p+1}) + \underline{T}_{k+1}(\sigma_{p+1}) + \sum_{j=k+2}^n \underline{T}_j(\sigma_p) \end{aligned} \quad (17)$$

Using (16), we note that  $\underline{T}_k(\sigma_p) = (a - B)^+$ ,  $\underline{T}_{k+1}(\sigma_p) = (A - b)^+$ ,  $\underline{T}_k(\sigma_{p+1}) = (a - b)^+$ , and  $\underline{T}_{k+1}(\sigma_{p+1}) = (A - B)^+$  if we set

$$\begin{aligned} a &= \frac{1}{C_i} \sum_{j=1}^k p_{i\pi_i(j)} c_{i\pi_i(j)}, & A &= \frac{1}{C_i} \sum_{j=1}^{k+1} p_{i\pi_i(j)} c_{i\pi_i(j)} \\ b &= d_{\sigma_p(k+1)}, & B &= d_{\sigma_p(k)} \end{aligned}$$

Note that  $a \leq A$ . Also,  $b \leq B$  since  $\sigma_p(k) > \sigma_p(k+1)$  and  $d_1 \leq \dots \leq d_n$ . From (17) we have

$$\underline{T}(\sigma_p) - \underline{T}(\sigma_{p+1}) = (a - B)^+ + (A - b)^+ - (a - b)^+ - (A - B)^+$$

It is straightforward to check that this quantity is always nonnegative when  $a \leq A$  and  $b \leq B$ . The theorem follows.

The bound of Lemma 2 can be written in terms of the variables  $x_{ik}$ :

$$\sum_{k=1}^n \underline{T}'_{ik} x_{ik}$$

where

$$\underline{T}'_{ik} \geq \frac{1}{C_i} \sum_{j=1}^k p_{i\pi_i(j)} c_{i\pi_i(j)} x_{i\pi_i(j)} - d_k, \quad k = 1, \dots, n$$

and  $\underline{T}'_{ik} \geq 0$ . We linearize the bound by writing it as

$$\sum_{k=1}^n \underline{T}_{ik} \tag{18}$$

where

$$\underline{T}_{ik} \geq \frac{1}{C_i} \sum_{j=1}^k p_{i\pi_i(j)} c_{i\pi_i(j)} x_{i\pi_i(j)} - d_k - (1 - x_{ik}) U_{ik}, \quad k = 1, \dots, n \tag{19}$$

and  $\underline{T}_{ik} \geq 0$ . The big-M term  $U_{ik}$  is given by

$$U_{ik} = \frac{1}{C_i} \sum_{j=1}^k p_{i\pi_i(j)} c_{i\pi_i(j)} - d_k$$

Note that although  $U_{ik}$  can be negative, the right-hand side of (19) is never positive when  $x_{ik} = 0$ . Finally, to obtain a relaxation of the subproblem, we sum (18) over all facilities and write

$$T \geq \sum_{i=1}^m \sum_{k=1}^n \underline{T}_{ik} \tag{20}$$

The relaxation (c) of the master problem now consists of (14), (20), and (19) for  $i = 1, \dots, m$ . The relaxation is valid only when tasks are indexed so that  $d_1 \leq \dots \leq d_n$ .

## 8 Problem Generation

Random instances were generated as follows. We set the number of facilities at 3, and the number of tasks at  $n = 10, 12, \dots, 24$ . The capacity limit was set to  $C_i = 10$  for each facility  $i$ . For each task  $j$ ,  $c_{ij}$  was assigned the same random value for all facilities  $i$  and drawn from a uniform distribution on  $[1, 10]$ . The processing time  $p_{ij}$  was drawn from a uniform distribution on  $[2, 20]$ ,  $[2, 25]$  and  $[2, 30]$  for facilities  $i = 1, 2, 3$ , respectively. For 22 or more tasks we used the intervals  $[5, 20]$ ,  $[5, 25]$  and  $[5, 30]$  since otherwise the minimum tardiness tends to be zero in the larger problems. The release dates were set to zero and the due date drawn from a uniform distribution on  $[\beta n/4, \beta n]$ . We used  $\beta = 20/9$ , partly since this was consistent with parameter settings used in earlier research, and partly because it leads to reasonable results (a few late tasks in most instances, and no late tasks in a few instances). No precedence constraints were used, which tends to make the scheduling portion of the problem more difficult.

## 9 Computational Results

We solved randomly generated problems with MILP (using CPLEX), CP (using the ILOG Scheduler), and the logic-based Benders method. All three methods were implemented with OPL Studio, using the OPL script language.

Table 1 shows computational results for minimizing the number of late tasks on three facilities using CP, MILP and the hybrid method. Since problem difficulty tends to increase with the minimum number of late tasks, the instances are ordered accordingly for each problem size. The problem instance identifier  $k$  appears in the last column. The instances are named  $ddn_j3mk$ , where  $n$  is the number of tasks and  $k$  the instance identifier. The instances are available at the web site [web.tepper.cmu.edu/jnh/planning.htm](http://web.tepper.cmu.edu/jnh/planning.htm).

On all but two problem instances the hybrid method is faster than MILP, which in turn is generally faster than CP. The advantage of the hybrid method becomes greater as the instances grow in size. The speedup is generally two or three orders of magnitude for instances with 16 or more tasks. The average speedup factor relative to MILP is 295 for these instances. This is almost certainly a substantial underestimate for the instances averaged, since the MILP solver was cut off after two hours. (The average omits instances in which the hybrid method was also cut off.) In addition MILP failed to solve 10 instances, while the hybrid method failed to solve only one instance.

Table 2 shows computational results for minimizing total tardiness. Again the hybrid method is almost always faster than MILP, which is faster than CP. The advantage of the hybrid approach is not as great as in the previous table, but the speedup factor is still significant on instances with 16 or more tasks. The average speedup factor on these instances is 25, which is again an underestimate for these instances. (The average omits instances for which the hybrid method was also cut off.)

The hybrid method failed to solve 6 of the 40 instances to optimality, only a modest improvement over the 10 that were intractable for MILP. However, when

**Table 1.** Computational results for minimizing the number of late tasks on three facilities. Computation is terminated after two hours (7200 seconds).

Tasks	Time (sec)			Hybrid/ MILP speedup	Best solution value found <sup>1</sup>		Instance
	CP	MILP	Hybrid		MILP	Hybrid	
10	0.09	0.48	0.05	9.6	<b>1</b>	<b>1</b>	1
	2.5	0.51	0.17	3.0	<b>1</b>	<b>1</b>	2
	0.28	0.46	0.27	1.7	<b>2</b>	<b>2</b>	5
	0.15	0.41	0.93	0.4	<b>3</b>	<b>3</b>	4
	1.7	3.9	3.0	1.3	<b>3</b>	<b>3</b>	3
12	0.01	0.73	0.07	10	<b>0</b>	<b>0</b>	1
	0.01	0.70	0.22	3.2	<b>0</b>	<b>0</b>	5
	0.02	0.64	0.06	11	<b>1</b>	<b>1</b>	3
	3.2	1.4	0.18	7.8	<b>1</b>	<b>1</b>	4
	1.6	1.7	0.34	5.0	<b>1</b>	<b>1</b>	2
14	1092	5.8	0.52	11	<b>1</b>	<b>1</b>	3
	382	8.0	0.69	12	<b>1</b>	<b>1</b>	2
	265	3.2	0.69	4.6	<b>2</b>	<b>2</b>	1
	85	2.6	1.3	2.0	<b>2</b>	<b>2</b>	5
	5228	1315	665	2.0	<b>3</b>	<b>3</b>	4
16	304	2.7	0.51	5.3	<b>0</b>	<b>0</b>	2
	? <sup>2</sup>	31	0.24	129	<b>1</b>	<b>1</b>	4
	310	22	0.41	54	<b>1</b>	<b>1</b>	5
	4925	29	2.7	11	<b>2</b>	<b>2</b>	3
	19	5.7	24	0.2	<b>4</b>	<b>4</b>	1
18	>7200	2.0	0.11	18	<b>0</b>	<b>0</b>	5
	? <sup>2</sup>	8.0	0.21	38	<b>1</b>	<b>1</b>	4
	>7200	867	8.5	102	<b>1</b>	<b>1</b>	2
	>7200	6.3	1.4	4.5	<b>2</b>	<b>2</b>	3
	>7200	577	3.4	170	<b>2</b>	<b>2</b>	1
20		97	0.37	262	<b>0</b>	<b>0</b>	1
		>7200	2.3	>3130	<b>1</b>	<b>1</b>	5
		219	5.0	44	<b>1</b>	<b>1</b>	2
		>7200	11	>655	<b>2</b>	<b>2</b>	3
		843	166	5.1	<b>3</b>	<b>3</b>	4
22		16	1.3	12	<b>0</b>	<b>0</b>	4
		>7200	3.7	>1946	<b>1</b>	<b>1</b>	1
		>7200	49	>147	<b>3</b>	<b>2</b>	5
		>7200	3453	>2.1	<b>5</b>	<b>2</b>	3
		>7200	>7200		<b>6</b>	<b>6</b>	2
24		25	0.8	31	<b>0</b>	<b>0</b>	3
		>7200	18	>400	<b>1</b>	<b>0</b>	5
		>7200	62	>116	<b>2</b>	<b>0</b>	4
		>7200	124	>58	<b>3</b>	<b>1</b>	1
		>7200	234	>31	<b>2</b>	<b>1</b>	2

<sup>1</sup>Values in boldface are proved optimal.

<sup>2</sup>Computation terminates with a segmentation fault.

the hybrid method failed to find provably optimal solutions, it obtained much better feasible solutions than obtained by MILP in the same two-hour period. In most cases these solutions were found very early in the solution process. Table 2 also shows the lower bounds obtained from the master problem, which in these instances are not very tight.

Table 3 illustrates the importance of relaxations in the hybrid approach, particularly when minimizing total tardiness. Lemmas 1 and 2 are clearly critical to the success of the hybrid method, especially when there are more than 16 tasks or so.

## 10 Conclusions

We find that integrating CP and MILP through a Benders scheme can substantially improve on the state of the art in planning and scheduling to minimize tardiness. The hybrid method is often two or three orders of magnitude faster than CP or MILP when minimizing the number of late tasks, and it solves significantly more problems. It is roughly an order of magnitude faster when minimizing total tardiness, and when it fails to solve the problem to optimality, it nonetheless finds a much better feasible solution in the same time period.

The problems become hard for all the methods examined when there are more than a few late tasks in the optimal solution. However, in such cases it is probably best to relax some of the time windows so as to reflect scheduling priorities, perhaps by postponing due dates for less critical tasks. This makes the problem easier to solve and yields a more meaningful compromise solution in practice.

## References

1. Benders, J. F., Partitioning procedures for solving mixed-variables programming problems, *Numerische Mathematik* **4** (1962) 238–252.
2. Cambazard, H., P.-E. Hladik, A.-M. Déplanche, N. Jussien, and Y. Trinquet, Decomposition and learning for a hard real time task allocation algorithm, in M. Wallace, ed., *Principles and Practice of Constraint Programming (CP 2004)*, *Lecture Notes in Computer Science* **3258**, Springer (2004).
3. Geoffrion, A. M., Generalized Benders decomposition, *Journal of Optimization Theory and Applications* **10** (1972) 237–260.
4. Hooker, J. N., *Logic-based Methods for Optimization: Combining Optimization and Constraint Satisfaction*, John Wiley & Sons (2000).
5. Hooker, J. N., A hybrid method for planning and scheduling, in M. Wallace, ed., *Principles and Practice of Constraint Programming (CP 2004)*, LNCS **3258**, 305–316.
6. Hooker, J. N., A hybrid method for planning and scheduling, *Constraints*, to appear.
7. Hooker, J. N. and G. Ottosson, Logic-based Benders decomposition, *Mathematical Programming* **96** (2003) 33–60.

**Table 2.** Computational results for minimum tardiness problems on three facilities. Computation is terminated after two hours (7200 seconds).

Tasks	Time (sec)			Hybrid/ MILP speedup	Best solution value found <sup>1</sup>		Benders lower bound <sup>2</sup>	Instance
	CP	MILP	Hybrid		MILP	Benders		
10	13	4.7	2.6	1.8	<b>10</b>	<b>10</b>		2
	1.1	6.4	1.6	4.0	<b>10</b>	<b>10</b>		1
	1.4	6.4	1.6	4.0	<b>16</b>	<b>16</b>		4
	4.6	32	4.1	7.8	<b>17</b>	<b>17</b>		5
	8.1	33	22	1.5	<b>24</b>	<b>24</b>		3
12	4.7	0.7	0.2	3.5	<b>0</b>	<b>0</b>		5
	14	0.6	0.1	6.0	<b>0</b>	<b>0</b>		1
	25	0.7	0.2	3.5	<b>1</b>	<b>1</b>		3
	19	15	2.4	6.3	<b>9</b>	<b>9</b>		4
	317	25	12	2.1	<b>15</b>	<b>15</b>		2
14	838	7.0	6.1	1.2	<b>1</b>	<b>1</b>		2
	7159	34	3.7	9.2	<b>2</b>	<b>2</b>		3
	1783	45	19	2.4	<b>15</b>	<b>15</b>		5
	> 7200	73	40	1.8	<b>19</b>	<b>19</b>		1
	> 7200	> 7200	3296	>2.2	26	<b>26</b>		4
16	> 7200	19	1.4	14	<b>0</b>	<b>0</b>		2
	> 7200	46	2.1	22	<b>0</b>	<b>0</b>		5
	> 7200	52	4.2	12	<b>4</b>	<b>4</b>		4
	> 7200	1105	156	7.1	<b>20</b>	<b>20</b>		3
	> 7200	3424	765	4.5	<b>31</b>	<b>31</b>		1
18		187	2.8	67	<b>0</b>	<b>0</b>		5
		15	5.3	2.8	<b>3</b>	<b>3</b>		4
		46	49	0.9	<b>5</b>	<b>5</b>		3
		256	47	5.5	<b>11</b>	<b>11</b>		1
		> 7200	1203	>6.0	14	<b>11</b>		2
20		105	18	5.8	<b>0</b>	<b>0</b>		1
		4141	23	180	<b>1</b>	<b>1</b>		5
		39	29	1.3	<b>4</b>	<b>4</b>		2
		1442	332	4.3	<b>8</b>	<b>8</b>		3
		> 7200	> 7200		75	37	9	4
22		6.3	19	0.3	<b>0</b>	<b>0</b>		4
		584	37	16	<b>2</b>	<b>2</b>		1
		> 7200	> 7200		120	40	7	3
		> 7200	> 7200		162	46	11	5
		> 7200	> 7200		375	141 <sup>3</sup>	34	2
24		10	324	0.03	<b>0</b>	<b>0</b>		3
		> 7200	94	>77	20	<b>0</b>		5
		> 7200	110	>65	57	<b>0</b>		4
		> 7200	> 7200		20	5	3	2
		> 7200	> 7200		25	7	1	1

<sup>1</sup>Values in boldface are proved optimal.

<sup>2</sup>When omitted, the lower bound is equal to the optimal value shown in the previous column.

<sup>3</sup>Best known solution is 128, obtained using a slightly weaker relaxation.

**Table 3.** Effect of relaxations on performance of the hybrid method. Computation time in seconds is shown.

Tasks	Minimizing late tasks:		Minimizing tardiness:		Instance
	with relaxation	without relaxation	with relaxation	without relaxation	
16	0.5	2.6	1.4	4.4	2
	0.4	1.5	2.1	6.5	5
	0.2	1.3	4.2	30	4
	2.7	4.2	156	199	3
	24	18	765	763	1
18	0.1	1.1	2.8	10	5
	0.2	0.7	5.3	17	4
	3.4	3.3	47	120	1
	1.4	15	49	354	3
	8.5	11	1203	5102	2
20	0.4	88	18	151	1
	2.3	9.7	23	1898	5
	5.0	63	29	55	2
	11	19	332	764	3
	166	226	>7200	>7200	4

8. Hooker, J. N., and Hong Yan, Logic circuit verification by Benders decomposition, in V. Saraswat and P. Van Hentenryck, eds., *Principles and Practice of Constraint Programming: The Newport Papers*, MIT Press (Cambridge, MA, 1995) 267–288.
9. Hooker, J. N., and Hong Yan, A relaxation for the cumulative constraint, in P. Van Hentenryck, ed., *Principles and Practice of Constraint Programming (CP2002)*, *Lecture Notes in Computer Science* **2470** (2002) 686–690.
10. Jain, V., and I. E. Grossmann, Algorithms for hybrid MILP/CP models for a class of optimization problems, *INFORMS Journal on Computing* **13** (2001) 258–276.
11. Thorsteinsson, E. S., Branch-and-Check: A hybrid framework integrating mixed integer programming and constraint logic programming, *Lecture Notes in Computer Science* **2239** (2001) 16–30.
12. Türkay, M., and I. E. Grossmann, Logic-based MINLP algorithms for the optimal synthesis of process networks, *Computers and Chemical Engineering* **20** (1996) 959–978.