

# Planning and Scheduling by Logic-Based Benders Decomposition

J. N. Hooker  
Carnegie Mellon University  
jh38@andrew.cmu.edu

July 2004

## Abstract

We combine mixed integer linear programming (MILP) and constraint programming (CP) to solve an important class of planning and scheduling problems. Tasks are allocated to facilities using MILP and scheduled using CP, and the two are linked via logic-based Benders decomposition. Tasks assigned to a facility may run in parallel subject to resource constraints (cumulative scheduling). We solve problems in which the objective is to minimize cost, makespan, or total tardiness. We obtain significant computational speedups, of several orders of magnitude for the first two objectives, relative to the state of the art in both MILP and CP. We also obtain better solutions and bounds for problems than cannot be solved to optimality.

We address a fundamental 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 deadlines. Some of the tasks assigned to a given facility may run in parallel if desired, provided the total resource consumption at any time remains within limits (“cumulative scheduling”). Precedence and other side constraints may be imposed. The objective is to minimize cost, makespan, or total tardiness.

Although problems of this sort arise often, they have proved to be quite difficult to solve. In the absence of an effective formal algorithm, a common-sense approach is often used in practice. The problem is decomposed into an assignment portion and a scheduling portion. Central managers assign tasks to facilities, whereupon operations managers develop detailed schedules for the tasks assigned to each facility. Since the schedule for one or more facilities may prove infeasible or otherwise unsuitable, the schedulers may telephone headquarters and ask for a different allocation of tasks. The give-and-take continues until a mutually satisfactory solution is reached. However, an informal process of this sort is time consuming, and even if carried to completion it offers no guarantee of obtaining a good solution.

One possibility that suggests itself is to apply Benders decomposition, which mimics the informal feedback loop in a rigorous fashion. The Benders master problem allocates tasks to facilities, and the subproblem separates into several independent scheduling problems. Benders cuts added to the master problem are a mathematical equivalent of telephone calls from the operations managers. Under typical conditions the Benders algorithm terminates with an optimal solution of the overall allocation and scheduling problem.

The classical Benders approach (Benders 1962, Geoffrion 1972) is inappropriate for this problem, however, because it requires that the subproblem be a continuous linear or nonlinear programming problem. Scheduling is a highly combinatorial problem that has no practical linear or nonlinear programming model. Fortunately, the idea of Benders decomposition can be extended to a *logic-based* form that accommodates an arbitrary subproblem, such as a discrete scheduling problem. This paper explores how a logic-based Benders method can be applied to planning and scheduling problems.

An additional advantage of logic-based Benders is that it permits one to exploit the relative strengths of mixed integer linear programming (MILP) and constraint programming (CP). MILP is a natural approach to the allocation task of the master problem, and CP provides an effective scheduling technology that can be applied to the subproblem. Logic-based Benders therefore unites MILP and CP in an integrated approach that obtains provably optimal solutions.

The planning and scheduling problem can be formulated entirely as an MILP problem, or entirely as a CP problem. However, these models are hard to solve. By linking MILP and CP, we obtained substantial speedups relative to the existing state of the art in both MILP (as represented by CPLEX) and CP (as represented by the ILOG Scheduler).

In particular, we obtained speedups of several orders of magnitude for minimum cost and minimum makespan problems. This allowed us to solve larger instances to optimality than could be solved previously. For minimum tardiness problems the computational advantage of a Benders method is about one order of magnitude. Although we were not able to solve significantly larger instances to optimality, we obtained much better solutions and better bounds on the optimal value, within a given amount of computation time, than provided by MILP and CP. One advantage of the Benders method is that it can be terminated at almost any point with a feasible solution and a lower bound on the optimal value (except in the case of the minimum cost problems). The bound improves steadily as the algorithm runs.

Unlike classical Benders, logic-based Benders provides no standard scheme for generating Benders cuts. Cuts must be devised for each problem class, and a major goal of this paper is to show that this can be done for a variety of objective functions. Specifically we develop Benders cuts for (a) minimum cost problems in which cost is a function of the assignment alone; (b) minimum makespan problems in which all tasks have the same release date; (c) minimum total tardiness problems in which all tasks have the same release date. We also develop relaxations of the subproblem for inclusion in the master problem, since this often

proves important for good performance. We performed computational tests on minimum cost and minimum makespan problems in which all tasks have the same time windows, and on tardiness problems in which all tasks have the same release date.

## 1 Previous Work

Classical Benders decomposition solves a problem by enumerating values of certain *primary* variables. For each set of values enumerated, it solves the subproblem that results from fixing the primary variables to these values. (In our case, the primary variables define the allocation of tasks to facilities.) Solution of the subproblem generates a Benders cut that the primary variables must satisfy in all subsequent solutions enumerated. The Benders cut is a linear inequality based on Lagrange multipliers obtained from a solution of the subproblem dual. The next set of values for the primary variables is obtained by solving the master problem, which contains all the Benders cuts so far generated. The process continues until the master problem and subproblem converge in value.

Logic-based Benders decomposition was introduced by Hooker and Yan (1995) in the context of logic circuit verification. The idea was formally developed by Hooker (2000) and applied to 0-1 programming by Hooker and Ottosson (2003). In logic-based Benders, the Benders cuts are obtained by solving the *inference dual* of the subproblem, of which the linear programming dual is a special case. The solution of the inference dual is a proof of optimality when the primary variables are fixed to their current values. The Benders cut is constructed by using this same proof to derive a valid bound on the optimal value when the primary variables take other values. Although logic-based Benders cuts can assume any form, in the present context they must be formulable as linear inequalities, since the master problem is an MILP.

Jain and Grossmann (2001) successfully applied logic-based Benders to minimum-cost planning and scheduling problems in which the subproblems are disjunctive scheduling problems in which tasks must be 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. A goal of the present paper is to show that a Benders approach can succeed for a variety of planning and scheduling problems in which the Benders cuts are less obvious. In particular, we (a) accommodate cumulative scheduling, and (b) develop Benders cuts for two additional objective functions for which the subproblem is an optimization problem.

In related work, Hooker (2000) observed that the master problem need only be solved once by a branching algorithm that accumulates Benders cuts as they are generated. Thorsteinsson (2001) showed that this approach, which he called *branch-and-check*, can result in substantially better performance on the Jain and Grossmann problems than standard logic-based Benders. We did not implement branch and check for this study because

it would require hand coding of a branch-and-cut algorithm for the master problem. But we obtained substantial speedups without it.

Classical Benders decomposition can also be useful in a CP context, as shown by Eremin and Wallace (2001).

## 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 deadline  $d_j$  (viewed as a due date in tardiness problems). 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

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

where  $J_{it} = \{j \mid s_j \leq t \leq s_j + p_{ij}\}$  is the set of tasks underway at time  $t$  in facility  $i$ .

An important special case is that in which the scheduling problem on facility  $i$  is an  $m$ -machine scheduling problem, for which  $C = m$  and each  $c_{ij} = 1$ . For a one-machine (“disjunctive”) scheduling problem we have  $C_i = 1$ .

Precedence and other side constraints may be added if desired. For instance, 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 three objective functions:

- *cost*, given by  $g(x, s) = \sum_j F_{x_j j}$ , where  $F_{ij}$  is the fixed cost of processing task  $j$  on facility  $i$ ;
- *makespan*, given by  $g(x, s) = \max_j \{s_j + p_{x_j j}\}$ ;
- *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.

When the tardiness objective is used, the deadlines become due dates and constraint (a) becomes  $r_j \leq s_j$  for all  $j$ .

### 3 Constraint Programming Formulation

The problem can be written in a constraint programming fashion using the “global” constraints *element* and *cumulative*. Given a vector  $v$  of values, the constraint

$$\text{element}(y, v, u) \tag{2}$$

sets  $u$  equal to the  $y$ th component of  $v$ . Thus one can implement an expression of the form  $v_y$  by replacing  $v_y$  with  $u$  and writing constraint (2).

If  $s = (s_1, \dots, s_n)$  and similarly for  $p$  and  $c$ , the constraint

$$\text{cumulative}(s, p, c, C) \tag{3}$$

requires that tasks be scheduled at times  $s_1, \dots, s_n$  so that the total rate of resource consumption at any given time never exceeds  $C$ . Thus  $\sum_{j \in J_t} c_j \leq C$  for all  $t$ , where  $J_t = \{j \mid s_j \leq t \leq s_j + p_{ij}\}$  is the set of tasks underway at time  $t$ . The cumulative constraint is widely used for resource-constrained scheduling problems in which tasks can run simultaneously.

Constraint programming solvers rely heavily on efficient “filtering” or domain reduction algorithms for these and other constraints. The filtering algorithms remove from the domain of each  $s_j$  some or all values that cannot be part of a solution satisfying the constraint.

The planning and scheduling problem now becomes

$$\begin{aligned} &\text{minimize } g(x, s) \\ &\text{subject to } r_j \leq s_j \leq d_j - u_j, \text{ all } j \\ &\quad \text{element}(x_j, (p_{1j}, \dots, p_{mj}), u_j), \text{ all } j \\ &\quad \text{cumulative}((s_j | x_j = i), (p_{ij} | x_j = i), (c_{ij} | x_j = i), C_i), \text{ all } i \end{aligned} \tag{4}$$

where  $g(x, s)$  is the desired objective function and  $(s_j | x_j = i)$  denotes the tuple of start times for tasks assigned to facility  $i$ . The first and second constraints enforce the time windows and the third enforces the resource constraint.

For purposes of computational testing we formulated (4) using the high-level modeling language of OPL Studio. The formulation for the minimum cost problem appears in Fig. 1. Formulations for the other objectives are similar. 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

```

[declarations]

scheduleHorizon = max(j in J) task[j].d;
DiscreteResource facility[i in I](C[i]);
AlternativeResources facilitySet(facility);
Activity scheduleTask[j in J];

minimize
    sum(j in J) cost[j]
subject to {
    forall(j in J)
        scheduleTask[j] requires(taskf[i,j].r) facilitySet;
    forall(j in J)
        forall(i in I)
            activityHasSelectedResource(scheduleTask[j],facilitySet,facility[i]) <=>
                scheduleTask[j].duration = taskf[i,j].p & cost[j] = taskf[i,j].F;

    forall(j in J) {
        scheduleTask[j].start >= task[j].r;
        scheduleTask[j].end <= task[j].d;
    };
};

search {
    assignAlternatives;
    setTimes;
};

```

Figure 1: OPL Studio code for the CP version of a minimum cost problem. Parameters  $r_j$ ,  $d_j$ ,  $p_{ij}$ , and  $F_{ij}$  are represented by `task[j].r`, `task[j].d`, `taskf[i,j].p`, and `taskf[i,j].F`, respectively. The resource limit  $C_i$  is `C[i]`. Index set  $J$  is  $\{1, \dots, n\}$ , and  $I$  is  $\{1, \dots, m\}$ .

discrete time  $t$  on facility  $i$ . The minimum cost formulation is

$$\begin{aligned}
 \min \quad & \sum_{ijt} F_{ij} x_{ijt} \\
 \text{subject to} \quad & \sum_{it} x_{ijt} = 1, \quad \text{all } j & (a) \\
 & \sum_j \sum_{t' \in T_{ijt}} c_{ij} x_{ijt'} \leq C_i, \quad \text{all } i, t & (b) \\
 & x_{ijt} = 0, \quad \text{all } j, t \text{ with } t \leq r_j \text{ or } t > d_j - p_{ij} & (c)
 \end{aligned} \tag{5}$$

where each  $x_{ijt}$  is a 0-1 variable.  $T_{ijt} = \{t' \mid t - p_{ij} < t' \leq t\}$  is the set of discrete times at which a task  $j$  in progress on facility  $i$  at time  $t$  might start processing. Constraint (a)

ensures that each task starts once on one facility, (b) enforces the resource limit, and (c) the time windows. The minimum makespan problem is

$$\begin{aligned}
& \min && M \\
& \text{subject to} && M \geq \sum_{it} (t + p_{ij}) x_{ijt}, \quad \text{all } j \\
& && \sum_{it} x_{ijt} = 1, \quad \text{all } j \\
& && \sum_j \sum_{t' \in T_{ijt}} c_{ij} x_{ijt'} \leq C_i, \quad \text{all } i, t \\
& && x_{ijt} = 0, \quad \text{all } j, t \text{ with } t \leq r_j \text{ or } t > d_j - p_{ij}
\end{aligned} \tag{6}$$

The minimum total tardiness problem is

$$\begin{aligned}
& \min && \sum_j T_j \\
& \text{subject to} && T_j \geq \sum_{it} (t + p_{ij}) x_{ijt} - d_j, \quad \text{all } j \\
& && \sum_{it} x_{ijt} = 1, \quad \text{all } j \\
& && \sum_j \sum_{t' \in T_{ijt}} c_{ij} x_{ijt'} \leq C_i, \quad \text{all } i, t \\
& && x_{ijt} = 0, \quad \text{all } j, t \text{ with } t \leq r_j \text{ or } t > d_j - p_{ij} \\
& && T_j \geq 0, \quad \text{all } j
\end{aligned} \tag{7}$$

where  $T_j$  is the tardiness of task  $j$ .

Models (5)-(7) can grow large due to time-indexed variables. We therefore investigated a smaller discrete event model suggested by Türkay and Grossmann (1996), which uses continuous time. Given  $n$  tasks, there are  $2n$  events. The model decides, for each event, which task it is associated with, whether it is the start or end of the task, and when the task starts. Let binary variable  $x_{ijk}$  be 1 when event  $k$  is the start of task  $j$  on facility  $i$ , and let  $y_{ijk}$  be 1 when event  $k$  is the end of task  $j$  on facility  $i$ . Continuous variable  $s_{ik}$  is the start time of event  $k$  on facility  $i$  when the event is assigned to facility  $i$ , and is zero otherwise. Continuous variable  $f_{ij}$  is the finish time of task  $j$  on facility  $i$  when task  $j$  is assigned to facility  $i$ , and is zero otherwise. The disaggregation of the start and finish time variables allows the model to impose resource constraints.

The minimum cost model is

$$\begin{aligned}
\min \quad & \sum_{ijk} F_{ij} x_{ijk} \\
\text{subject to} \quad & \sum_{ik} x_{ijk} = 1, \quad \sum_{ik} y_{ijk} = 1, \quad \text{all } j \quad (a) \\
& \sum_{ij} x_{ijk} + y_{ijk} = 1, \quad \text{all } k \quad (b) \\
& \sum_k x_{ijk} = \sum_k y_{ijk}, \quad \text{all } i, j \quad (c) \\
& s_{i,k-1} \leq s_{ik}, \quad \text{all } i, k \quad (d) \\
& r_j x_{ijk} \leq s_{ik}, \quad f_{ij} \leq d_j, \quad \text{all } i, j, k \quad (e) \\
& s_{ik} + p_{ij} x_{ik} - M(1 - x_{ijk}) \leq f_{ij} \\
& \qquad \qquad \qquad \leq s_{ik} + p_{ik} x_{ijk} + M(1 - x_{ijk}), \quad \text{all } i, j, k \quad (f) \\
& s_{ik} - M(1 - y_{ijk}) \leq f_{ij} \leq s_{ik} + M(1 - y_{ijk}), \quad \text{all } i, j, k \quad (g) \\
& R_{ik} \leq C_i, \quad \text{all } i, k \quad (h) \\
& R_{i1} = R_{i1}^s, \quad R_{ik}^s = \sum_j c_{ij} x_{ijk}, \quad R_{ik}^f = \sum_j c_{ij} y_{ijk}, \quad \text{all } i, k \quad (i) \\
& R_{ik}^s + R_{i,k-1} - R_{ik}^f = R_{ik}, \quad \text{all } i, k \quad (j) \\
& x_{ijk}, y_{ijk} \in \{0, 1\}, \quad R_{ik}^s, R_{ik}^f \geq 0, \quad \text{all } i, k
\end{aligned}$$

Constraints (a), (b) and (c) ensure that each task is assigned to one facility, starts and ends once, and starts and ends on the same facility. Constraint (d) requires events  $1, \dots, 2n$  to occur in chronological order. Constraint (e) enforces the time windows. Constraints (f) and (g) define the finish time variables in terms of the start time variables.

The rest of the model encodes the resource constraints. Variable  $R_{ik}$  is the total resource consumption of tasks running on facility  $i$  at the time that event  $k$  occurs. Constraint (h) imposes the upper bound on resource consumption at the time of each event, which is sufficient to impose the bound at all times. Constraint (i) computes the resource consumption  $R_{ik}^s$  of any task starting at the time of event  $k$  and the consumption  $R_{ik}^f$  of any task finishing at the time of event  $k$ . Constraint (j) maintains  $R_{ik}$  at the time of each event in the fashion of an inventory balance equation.

Although it is smaller than (5), the discrete event model proved to be much harder to solve. We therefore omitted it from the computational results given below.



## 5 Logic-based Benders Decomposition

Logic-based Benders decomposition applies to problems of the form

$$\begin{aligned} & \text{minimize} && f(x, y) \\ & \text{subject to} && C(x, y) \\ & && x \in D_x, y \in D_y \end{aligned} \tag{8}$$

where  $C(x, y)$  is a set of constraints containing variables  $x, y$ .  $D_x$  and  $D_y$  denote the domains of  $x$  and  $y$ , respectively. When  $x$  is fixed to a given value  $\bar{x} \in D_x$ , the following *subproblem* results:

$$\begin{aligned} & \text{minimize} && f(\bar{x}, y) \\ & \text{subject to} && C(\bar{x}, y) \\ & && y \in D_y \end{aligned} \tag{9}$$

Here  $C(\bar{x}, y)$  is the constraint that results from fixing  $x = \bar{x}$  in  $C(x, y)$ .

The *inference dual* of (9) is the problem of inferring the tightest possible lower bound on  $f(\bar{x}, y)$  from  $C(\bar{x}, y)$ . It can be written

$$\begin{aligned} & \text{maximize} && v \\ & \text{subject to} && C(\bar{x}, y) \implies f(\bar{x}, y) \geq v \\ & && v \in R \end{aligned} \tag{10}$$

where  $\implies$  means “implies.” The inference dual becomes the linear programming dual when (9) is a linear programming problem; see Hooker (2000) for details.

The solution of the dual can be viewed as a derivation of the tightest possible bound  $\hat{v}$  on  $f(x, y)$  when  $x = \bar{x}$ . For purposes of Benders decomposition, we wish to derive not only a bound when  $x = \bar{x}$  but a bounding function  $B_{\bar{x}}(x)$  with two properties:

B1.  $B_{\bar{x}}(x)$  provides a valid lower bound on  $f(x, y)$  for any given  $x \in D_x$ . That is,  $f(x, y) \geq B_{\bar{x}}(x)$  for any feasible  $(x, y)$  in (8).

B2. In particular,  $B_{\bar{x}}(\bar{x}) = \hat{v}$ .

If  $z$  is the objective function value of (8), the valid inequality  $z \geq B_{\bar{x}}(x)$  is a *Benders cut*.

In iteration  $H$  of the Benders algorithm, we solve a *master problem* whose constraints are the Benders cuts so far generated:

$$\begin{aligned} & \min && z \\ & \text{subject to} && z \geq B_{x^h}(x), h = 1, \dots, H - 1 \\ & && z \in R, x \in D_x \end{aligned} \tag{11}$$

Here  $x^1, \dots, x^{H-1}$  are the solutions of the previous  $H - 1$  master problems. Then the solution  $\bar{x}$  of (11) defines the next subproblem (9).

If we let  $v_1^*, \dots, v_{H-1}^*$  denote the optimal values of the previous  $H - 1$  subproblems, the algorithm continues until the optimal value  $z_H^*$  of the master problem equals  $v^* = \min\{v_1^*, \dots, v_{H-1}^*\}$ . At any point in the algorithm,  $z_H^*$  and  $v^*$  provide lower and upper bounds on the optimal value of the problem.

The following is shown by Hooker (2000). It is convenient to regard (8) or (9) as having an infinite optimal value if it is infeasible.

**Theorem 1** *If the bounding function  $B_{\bar{x}}(x)$  satisfies (B1) and (B2) in each iteration of the Benders algorithm, and  $D_y$  is finite, the Benders algorithm converges to the optimal value of (8) after finitely many steps.*

In the planning and scheduling problem (1), variables  $x, s$  correspond respectively to  $x, y$  in (8). Any assignment  $\bar{x}$  of tasks to facilities creates the subproblem:

$$\begin{aligned} \min \quad & g(\bar{x}, s) \\ \text{subject to} \quad & \text{cumulative}((s_j | \bar{x}_j = i), (p_{ij} | \bar{x}_j = i), (c_{ij} | \bar{x}_j = i), C_i), \text{ all } i \\ & r_j \leq s_j \leq d_j - p_{\bar{x}_j j}, \text{ all } j \end{aligned} \quad (12)$$

which decomposes into a separate scheduling problem for each facility  $i$ . After solving the subproblem with a finite-domain CP solver, we generate a Benders cut that becomes part of the master problem (11), which is formulated as an MILP. The domain  $D_s$  can be regarded as finite, as required in Theorem 1, since the CP solver assumes a finite set of start times (such as integers). Release times and deadlines are typically given as integers.

The bounding function  $B_{\bar{x}}(x)$  is generally obtained by examining the type of reasoning that led to a bound for  $x = \bar{x}$  and extending this reasoning to obtain a bound for general  $x$ . In the present context, however, only the primal solution (the schedule itself) is available from the commercial CP solver. We therefore restrict ourselves in this study to Benders cuts that can be formulated without internal information from the CP solver.

## 6 Minimizing Cost

The cost objective presents the simplest case, since cost can be computed solely in terms of master problem variables. The subproblem (12) takes the value  $\sum_j F_{x_j j}$  when it is feasible and  $\infty$  when infeasible. Let  $J_{hi} = \{j \mid x_j^h = i\}$  be the set of tasks assigned to facility  $i$  in iteration  $h$ . If there is no feasible schedule for facility  $i$ , the most obvious Benders cut simply rules out assigning this same set  $J_{hi}$  of tasks to facility  $i$ . In this case the bounding

function takes the form

$$B_{x^h}(x) = \begin{cases} \infty & \text{if the scheduling subproblem on some facility } i \text{ is infeasible} \\ & \text{and } x_j^h = x_j \text{ for all } j \in J_{hi} \\ \sum_j F_{x_j j} & \text{otherwise} \end{cases}$$

Although we do not know the reasoning that the CP solver used to prove that no schedule is feasible for  $J_{hi}$ , we know that this reasoning, whatever it is, proves that no schedule is feasible for any superset of  $J_{hi}$ . This is basis for the bounding function.

The function  $B_{x^h}$  satisfies condition (B1) because the cost  $\sum_j F_{x_j j}$  is equal to  $B_{x^h}(x)$  for any feasible  $(x, s)$  and is infinite for any infeasible  $(x, s)$ . It satisfies (B2) because  $B_{x^h}(x^h)$  is the optimal value of the subproblem.

The master problem (10), written as a 0-1 programming problem, becomes

$$\begin{aligned} & \text{minimize} && \sum_{ij} F_{ij} x_{ij} \\ & \text{subject to} && \sum_i x_{ij} = 1, \text{ all } j && (a) \\ & && \sum_{j \in J_{hi}} (1 - x_{ij}) \geq 1, \text{ all } i, h = 1, \dots, H-1 && (b) \\ & && \text{relaxation of subproblem} && (c) \end{aligned} \tag{13}$$

where  $x_{ij} \in \{0, 1\}$ , and where constraints (b) are the Benders cuts.

Experience shows that it is important to include a relaxation of the subproblem within the master problem. A straightforward relaxation can be obtained as follows. For any two times  $t_1, t_2$ , let  $J(t_1, t_2)$  be the set of tasks  $j$  whose time windows lie between  $t_1$  and  $t_2$ ; that is,  $[r_j, d_j] \subset [t_1, t_2]$ . If the tasks  $j \in J \subset J(t_1, t_2)$  are assigned to the same facility  $i$ , then clearly the ‘‘area’’  $\sum_{j \in J} p_{ij} c_{ij}$  of these tasks can be at most  $C_i(t_2 - t_1)$  if they are to be scheduled in the time interval  $[t_1, t_2]$ . This yields the valid inequality

$$\frac{1}{C_i} \sum_{j \in J(t_1, t_2)} p_{ij} c_{ij} x_{ij} \leq t_2 - t_1, \tag{14}$$

which we refer to as inequality  $R^i(t_1, t_2)$ . If we let  $\bar{r}_1, \dots, \bar{r}_{n_r}$  be the distinct elements of  $\{r_1, \dots, r_n\}$  in increasing order, and similarly for  $\bar{d}_1, \dots, \bar{d}_{n_d}$ , we have a relaxation consisting of the inequalities

$$R^i(\bar{r}_j, \bar{d}_k), \quad j = 1, \dots, n_r, \quad k = 1, \dots, n_d \tag{15}$$

for each facility  $i$ . These inequalities serve as the relaxation (c) in (13). There are doubtless stronger relaxations, but we require a linear relaxation that is expressed in terms of the variables  $x_{ij}$ .

Let  $\mathcal{R}_i = \emptyset$ .  
For  $j = 1, \dots, p$ :  
Set  $k' = 0$ .  
For  $k = 1, \dots, q$ :  
If  $r_k \geq r_j$  and  $T^i(\bar{r}_j, \bar{d}_k) < T^i(\bar{r}_j, \bar{d}_{k'})$  then  
Remove from  $\mathcal{R}_i$  all  $R^i(\bar{r}_{j'}, \bar{d}_k)$  for which  $T^i(\bar{r}_j, \bar{d}_k) \geq T^i(\bar{r}_{j'}, \bar{d}_k)$ .  
Add  $R^i(\bar{r}_j, \bar{d}_k)$  to  $\mathcal{R}_i$  and set  $k' = k$ .

Figure 2:  $O(n^3)$  algorithm for generating an inequality set  $\mathcal{R}_i$  that relaxes the time window constraints for facility  $i$ . By convention  $\bar{d}_0 = -\infty$ .

Let  $\mathcal{R}_i = \emptyset$ , and set  $j = 0$ .  
For  $k = 1, \dots, p_d$ :  
If  $T^i(r_0, \bar{d}_k) > T^i(r_0, \bar{d}_j)$  then add  $R^i(r_0, \bar{d}_k)$  to  $\mathcal{R}_i$  and set  $j = k$ .

Figure 3:  $O(n)$  algorithm for generating an inequality set  $\mathcal{R}_i$  that relaxes the time window constraints for facility  $i$ , where  $r_0 = \bar{d}_0 = \min_j\{r_j\}$  and  $T^i(r_0, r_0) = 0$ .

Many of these inequalities may be redundant of the others, and if desired they can be omitted from the relaxation. Let

$$T^i(t_1, t_2) = \frac{1}{C_i} \sum_{j \in J(t_1, t_2)} p_{ij} c_{ij} - t_2 + t_1$$

measure the “tightness” of  $R^i(t_1, t_2)$ . It is easily verified that  $R^i(t_1, t_2)$  dominates  $R^i(u_1, u_2)$  whenever  $[t_1, t_2] \subset [u_1, u_2]$  and  $T^i(t_1, t_2) \geq T^i(u_1, u_2)$ . Dominated inequalities can now be removed.

A set of undominated inequalities can be generated for each facility using the algorithm of Fig. 2. It has  $O(n^3)$  complexity in the worst case, since it is possible that none of the inequalities are eliminated. This occurs, for instance, when each  $r_j = j - 1$ ,  $d_j = j$ , and  $p_{ij} = 2$ . However, the algorithm need only be run once as a preprocessing routine.

In practice the relaxation can be simplified by supposing that the release times are all  $r_0 = \min_j\{r_j\}$ . Then the relaxation (c) in (13) consists of

$$R^i(r_0, \bar{d}_k), \quad k = 1, \dots, n_d$$

for each facility  $i$ . The redundant inequalities can be eliminated running the simple  $O(n)$  algorithm of Fig. 3 for each facility. Similarly, one can suppose that the deadlines are all  $d_0 = \max_j\{d_j\}$  and use the inequalities  $R^i(\bar{r}_j, d_0)$ .

Precedence constraints may be accommodated if they apply only to tasks assigned to the same facility. Thus if tasks  $j$  and  $k$  must be assigned to the same facility, and  $j$  must

precede  $k$ , we add the constraint  $x_j = x_k$  to the master problem (13) and  $t_j + p_{ij} \leq t_j$  for all  $i$  to the subproblem (12).

## 7 Minimizing Makespan

This case is less straightforward because the subproblem is an optimization problem. However, there are relatively simple linear Benders cuts when all tasks have the same release date, and they simplify further when all deadlines are the same. We also use a linear subproblem relaxation that is valid for any set of time windows.

The Benders cuts are based on the following fact:

**Lemma 2** *Consider a minimum makespan problem  $P$  in which tasks  $1, \dots, n$  with release time  $0$  and deadlines  $d_1, \dots, d_n$  are to be scheduled on a single facility  $i$ . Let  $M^*$  be the minimum makespan for  $P$ , and  $\hat{M}$  the minimum makespan for the problem  $\hat{P}$  that is identical to  $P$  except that tasks  $1, \dots, s$  are removed. Then*

$$M^* - \hat{M} \leq \Delta + \max_{j \leq s} \{d_j\} - \min_{j \leq s} \{d_j\} \quad (16)$$

where  $\Delta = \sum_{j=1}^s p_{ij}$ . In particular, when all the deadlines are the same,  $M^* - \hat{M} \leq \Delta$ .

*Proof.* Consider any optimal solution of  $\hat{P}$  and extend it to a solution  $S$  of  $P$  by scheduling tasks  $1, \dots, s$  sequentially after  $\hat{M}$ . That is, for  $k = 1, \dots, s$  let task  $k$  start at time  $\hat{M} + \sum_{j=1}^{k-1} p_{ij}$ . The makespan of  $S$  is  $\hat{M} + \Delta$ . If  $\hat{M} + \Delta \leq \min_{j \leq s} \{d_j\}$ , then  $S$  is clearly feasible for  $P$ , so that  $M^* \leq \hat{M} + \Delta$  and the lemma follows. Now suppose  $\hat{M} + \Delta > \min_{j \leq s} \{d_j\}$ . This implies

$$\hat{M} + \Delta + \max_{j \leq s} \{d_j\} - \min_{j \leq s} \{d_j\} > \max_{j \leq s} \{d_j\} \quad (17)$$

Since  $M^* \leq \max_{j \leq s} \{d_j\}$ , (17) implies (16), and again the lemma follows.

The bound  $M^* - \hat{M} \leq \Delta$  need not hold when the deadlines differ. Consider for example an instance with three tasks where  $(r_1, r_2, r_3) = (0, 0, 0)$ ,  $(d_1, d_2, d_3) = (2, 1, \infty)$ ,  $(p_{i1}, p_{i2}, p_{i3}) = (1, 1, 2)$ , and  $(c_{i1}, c_{i2}, c_{i3}) = (2, 1, 1)$ . Then if  $s = 1$ , we have  $M^* - \hat{M} = 4 - 2 > \Delta = p_{i1} = 1$ .

Now consider any given iteration of the Benders algorithm, and suppose for the moment that all the time windows are identical. Let  $J_{hi}$  be the set of tasks assigned to facility  $i$  in a previous iteration  $h$ , and  $M_{hi}^*$  the corresponding minimum makespan incurred by facility  $i$ . The solution of the current master problem removes task  $j \in J_{hi}$  from facility  $i$  when

$x_{ij} = 0$ . Thus by Lemma 2, the resulting minimum makespan for facility  $i$  is reduced by at most

$$\sum_{j \in J_{hi}} p_{ij}(1 - x_{ij}) \quad (18)$$

and we therefore have a lower bound on the makespan of facility  $i$ :

$$M_{hi}^* - \sum_{j \in J_{hi}} p_{ij}(1 - x_{ij})$$

This yields a bounding function

$$B_{x^h}(x) = \max_i \left\{ M_{hi}^* - \sum_{j \in J_{hi}} p_{ij}(1 - x_{ij}) \right\}$$

that bounds the overall makespan and therefore satisfies (B1). The function satisfies (B2) as well because  $B_{x^h}(x^h) = M_{hi}^*$ , which is the optimal subproblem value.

We can now write Benders cuts (b) in the master problem:

$$\begin{aligned} & \text{minimize } M \\ & \text{subject to } \sum_i x_{ij} = 1, \text{ all } j \quad (a) \\ & M \geq M_{hi}^* - \sum_{j \in J_{hi}} (1 - x_{ij})p_{ij}, \text{ all } i, h = 1, \dots, H - 1 \quad (b) \quad (19) \\ & M \geq \frac{1}{C_i} \sum_j c_{ij}p_{ij}x_{ij}, \text{ all } i \quad (c) \end{aligned}$$

where  $x_{ij} \in \{0, 1\}$ . The relaxation (c) is similar to that for the minimum cost problem.

If the deadlines differ (and the release times still the same), the minimum makespan for facility  $i$  is at least

$$M_{hi}^* - \left( \sum_{j \in J_{hi}} p_{ij}(1 - x_{ij}) + \max_{j \in J_{hi}} \{d_j\} - \min_{j \in J_{hi}} \{d_j\} \right)$$

if one or more tasks are removed, and is  $M_{hi}^*$  otherwise. The bounding function that results can be linearized to obtain the Benders cuts

$$\left. \begin{aligned} M &\geq M_{hi}^* - \sum_{j \in J_{hi}} (1 - x_{ij})p_{ij} - w_{hi} \\ w_{hi} &\leq \left( \max_{j \in J_{hi}} \{d_j\} - \min_{j \in J_{hi}} \{d_j\} \right) \sum_{j \in J_{hi}} (1 - x_{ij}) \\ w_{hi} &\leq \max_{j \in J_{hi}} \{d_j\} - \min_{j \in J_{hi}} \{d_j\} \end{aligned} \right\} \text{ all } i, h = 1, \dots, H - 1 \quad (20)$$

Thus (20) replaces (19b) when the deadlines differ and all release times are equal. For purposes of computational testing, however, we focus on the case in which all time windows are the same, since this seems to be the more important case, and it permits simpler and stronger Benders cuts.

## 8 Minimizing Total Tardiness

For this problem we again assume that all release times are zero. We consider two approaches to formulating Benders cuts, the first of which constructs cuts in a manner similar to those obtained for the minimum makespan problem. The following lemma provides the basis for these cuts:

**Lemma 3** *Consider a minimum total tardiness problem  $P$  in which tasks  $1, \dots, n$  have release time 0, due dates  $d_1, \dots, d_n$ , and are to be scheduled on a single facility  $i$ . Let  $T^*$  be the minimum tardiness for this problem, and  $\hat{T}$  the minimum tardiness for the problem  $\hat{P}$  that is identical to  $P$  except that tasks  $1, \dots, s$  are removed. Then*

$$T^* - \hat{T} \leq \sum_{k=1}^s \left( \sum_{j=1}^n p_{ij} - d_k \right)^+ \quad (21)$$

*Proof.* Consider any optimal solution  $\hat{S}$  of  $\hat{P}$ . We may assume that the makespan of  $\hat{S}$  is at most  $M = \sum_{j=s+1}^n p_{ij}$ , since if it is greater than  $M$ , we can move at least one task to an earlier time without reducing the total tardiness of  $\hat{S}$ . To obtain a feasible solution  $S$  for  $P$ , schedule tasks  $1, \dots, s$  sequentially after  $M$ . That is, for  $k = 1, \dots, s$  let task  $k$  start at time  $M + \sum_{j=1}^{k-1} p_{ij}$ . The tardiness of task  $k$  in  $S$  is at most

$$\left( M + \sum_{j=1}^s p_{ij} - d_k \right)^+ = \left( \sum_{j=1}^n p_{ij} - d_k \right)^+$$

The total tardiness of  $S$  is therefore at most

$$\hat{T} + \sum_{k=1}^s \left( \sum_{j=1}^n p_{ij} - d_k \right)^+$$

from which (21) follows.

The proof uses the fact that the makespan  $\hat{M}$  of  $\hat{S}$  can be assumed to be at most  $M$ . A generally tighter bound could be obtained by assuming  $\hat{M}$  is at most the makepan  $M^*$

of  $P$ 's minimum tardiness solution. Unfortunately, this is not a valid assumption, since in exceptional cases  $\hat{M} > M^*$  even though  $\hat{P}$  has fewer tasks than  $P$ . Suppose for example that  $n = 4$  and  $s = 1$ , with  $r = (0, 0, 0, 0)$ ,  $d = (5, 3, 3, 6)$ ,  $p = (1, 2, 2, 4)$ ,  $c = (2, 1, 1, 1)$ , and  $C = 2$ . An optimal solution of  $P$  puts  $t = (4, 0, 2, 0)$  with tardiness  $T^* = 1$  and makespan  $M^* = 5$ , but the only optimal solution of the smaller problem  $\hat{P}$  puts  $(t_2, t_3, t_4) = (0, 0, 2)$  with tardiness  $\hat{T} = 0$  and makespan  $\hat{M} = 6 > M^*$ .

Lemma 3 is applied as follows. As before let  $J_{hi}$  be the set of tasks assigned to facility  $i$  in a previous iteration  $h$ , and  $T_{hi}^*$  the corresponding minimum tardiness incurred by facility  $i$ . Since task  $j$  is removed from facility  $i$  when  $x_{ij} = 0$ , the resulting minimum tardiness on facility  $i$  is at least

$$T_{hi}^* - \sum_{k \in J_{hi}} \left( \sum_{j \in J_{hi}} p_{ij} - d_k \right)^+ (1 - x_{ik})$$

This yields the Benders cuts (b) in the master problem

$$\begin{aligned} & \text{minimize } T \\ & \text{subject to } \sum_i x_{ij} = 1, \text{ all } j \end{aligned} \tag{a}$$

$$T \geq T_{hi}^* - \sum_{k \in J_{hi}} \left( \sum_{j \in J_{hi}} p_{ij} - d_k \right)^+ (1 - x_{ik}), \text{ all } i, h = 1, \dots, H - 1 \tag{b}$$

$$\text{relaxation of subproblem} \tag{c}$$

(22)

Unfortunately these cuts are weak and do not perform well in practice, due to the lack of a tighter bound on  $\hat{M}$ . It appears that stronger cuts require some internal information from the CP solver.

A second approach to formulating Benders cuts acquires such internal information indirectly. One can re-solve the subproblems a few times with some tasks removed in order to observe how much the tardiness actually decreases. This incurs minimal computational cost, since the MILP master problem consumes the bulk of the computational time in the tardiness instances solved for this study.

Several such schemes are possible, but we found the following to be effective. For any subset  $J$  of the tasks  $J_{hi}$  assigned to facility  $i$ , let  $M_{hi}^*(J)$  be the minimum makespan that results when the tasks in  $J$  are removed, so that  $M_{hi}^*(\emptyset) = M_{hi}^*$ . Let  $J_{hi}^0$  be the set of tasks in  $J_{hi}$  that can be removed, one at a time, without reducing the minimum makespan. That is,

$$J_{hi}^0 = \{j \in J_{hi} \mid M_{hi}^*(\{j\}) = M_{hi}^*\}$$

Finally, let  $\Delta_{hi}$  be the reduction in minimum makespan that results from removing all the tasks in  $J_{hi}^0$  at once, so that  $\Delta_{hi} = M_{hi}^*(J_{hi}^0)$ . Thus any or all tasks in  $J_{hi}^0$  can be removed



from facility  $i$  without reducing the minimum makespan more than  $\Delta_{hi}$ . This yields a bounding function  $B_{x^h}(x) = \min_i B_{x^h}^i(x)$ , where

$$B_{x^h}^i(x) = \begin{cases} T_{hi}^* - \Delta_{hi} & \text{if } J_{hi}^0 \supset \{j \in J_{hi} \mid x_j \neq i\} \neq \emptyset & (a) \\ T_{hi}^* & \text{if } \{j \in J_{hi} \mid x_j \neq i\} = \emptyset & (b) \\ 0 & \text{otherwise} & (c) \end{cases} \quad (23)$$

The function  $B_{x^h}(x)$  satisfies (B1) by construction. Case (b) of (23) ensures that it satisfies (B2) as well.

The master problem MILP now becomes

$$\begin{aligned} & \text{minimize } T \\ & \text{subject to } \sum_i x_{ij} = 1, \text{ all } j & (a) \\ & T \geq T_{hi}^* - \Delta_{hi} - U \sum_{j \in J_{hi} \setminus J_{hi}^0} (1 - x_{ij}), \text{ all } i, h = 1, \dots, H-1 & (b1) \\ & T \geq T_{hi}^* - U \sum_{j \in J_{hi}} (1 - x_{ij}), \text{ all } i, h = 1, \dots, H-1 & (b2) \\ & \text{relaxation of subproblem} & (c) \end{aligned} \quad (24)$$

where the “big-M” coefficient  $U$  is a large number. The cut (b2) is redundant and can be eliminated for a given  $h, i$  when  $\Delta_{hi} = 0$ . This in fact substantially reduces the size of master problem, since computational testing suggests that  $\Delta_{hi}$  is usually zero. The computational work required in each iteration to derive the cuts consists of solving  $n$  scheduling subproblems to compute  $M_{hi}^*(\{j\})$  for each  $i$  and each  $j \in J_{hi}$ , plus  $m$  scheduling subproblems to compute  $M_{hi}^*(J_{hi}^0)$  for each  $i$ , where  $n$  is the number of tasks and  $m$  the number of facilities.

We employed two relaxations for the master problem (24), since they seem to have somewhat complementary strengths. The first and simpler relaxation is again based on the “areas” required by tasks. Recall that  $J(t_1, t_2)$  is the set of jobs with time windows between  $t_1$  and  $t_2$ .

**Lemma 4** *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  $\hat{T} = \frac{1}{C_i} \sum_{j \in J(0, d_k)} p_{ij} c_{ij}$ . Since the last task has due date no later than  $d_k$ , its tardiness is no less than  $(\hat{T} - d_k)^+$ . Thus total tardiness is no less than  $(\hat{T} - d_k)^+$ .

This gives rise to a relaxation consisting of

$$T \geq \frac{1}{C_i} \sum_{j \in J(0, d_k)} p_{ij} c_{ij} x_{ij} - d_k, \quad \text{all } i, k \quad (25)$$

and  $T \geq 0$ .

A second relaxation of the subproblem 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 5** *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

$$T = \sum_{k=1}^n (t_k + p_{ik} - d_k)^+ \quad (26)$$

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

$$\underline{T}_k(\sigma) = \left( \frac{1}{C_i} \sum_{j=1}^k p_{i\pi_i(j)} c_{i\pi_i(j)} - d_{\sigma(k)} \right)^+ \quad (27)$$

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 (26) 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 (27) 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

$$\begin{aligned} \underline{T}(\sigma_p) &= \sum_{j=1}^{k-1} \underline{T}_j(\sigma_p) + \underline{T}_k(\sigma_p) + \underline{T}_{k+1}(\sigma_p) + \sum_{j=k+2}^n \underline{T}_j(\sigma_p) \\ \underline{T}(\sigma_{p+1}) &= \sum_{j=1}^{k-1} \underline{T}_j(\sigma_p) + \underline{T}_k(\sigma_{p+1}) + \underline{T}_{k+1}(\sigma_{p+1}) + \sum_{j=k+2}^n \underline{T}_j(\sigma_p) \end{aligned} \tag{28}$$

Using (27), 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 (28) 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 Theorem 5 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

$$\sum_{k=1}^n \underline{T}_{ik} \tag{29}$$

where

$$\underline{T}_{ik} \geq \frac{1}{C_i} \sum_{j=1}^k p_{i\pi_i(j)} c_{i\pi_i(j)} x_{i\pi_i(j)} - d_k - (1 - x_{ik}) U_{ik}, \quad k = 1, \dots, n \tag{30}$$

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 (30) is never positive when  $x_{ik} = 0$ . Finally, to obtain a relaxation of the subproblem, we sum (29) over all facilities and write

$$T \geq \sum_{i=1}^m \sum_{k=1}^n \underline{T}_{ik} \tag{31}$$

The relaxation (c) of the master problem now consists of two relaxations, the first of which is (25), and the second (31) together with (30) for  $i = 1, \dots, m$ . The second relaxation is valid only when tasks are indexed so that  $d_1 \leq \dots \leq d_n$ .

## 9 Problem Generation

Random instances were generated as follows for minimum cost and minimum makespan problems. 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  $[1, 10]$ . For instances with  $n$  tasks and  $m$  facilities, the processing time  $p_{ij}$  of each task  $j$  on facility  $i$  was drawn from a uniform distribution on  $[i, 10i]$ . Thus facility 1 tends to run about  $i$  times faster than facility  $i$  for  $i = 1, \dots, m$ . Since the average of  $10i$  over  $m$  facilities is  $5(m+1)$ , the total processing time of all tasks is roughly proportional to  $5n(m+1)$ , or about  $5n(m+1)/m$  per facility. The release dates were set to zero, and the deadline for every task was set to  $5\alpha n(m+1)/m$  (rounded to the nearest integer). We used  $\alpha = 1/3$  for the minimum cost and minimum makespan problem, which results in a deadline that

is loose enough to permit feasible solutions but tight enough so that tasks are reasonably well distributed over the facilities in minimum cost solutions. In minimum cost problems, the cost  $F_{ij}$  is drawn from a uniform distribution on  $[2(m - i + 1), 20(m - i + 1)]$ , so that faster facilities tend to be more expensive.

No precedence constraints were used, which tends to make the scheduling portion of the problem more difficult. The wide variation in facility speeds tends to put the Benders approach at a disadvantage, since the faster facilities receive more tasks and require disproportionately more computation to solve the scheduling problem.

For minimum tardiness problems we set the release dates to zero and drew the due dates from a uniform distribution on the interval  $[(5/4)\alpha n(m+1)/m, 5\alpha n(m+1)/m]$ , again using  $\alpha = 1/3$ . We found that when the processing time for facility  $i$  is drawn from  $[i, 10i]$ , the wide range of facility speeds often results in solutions in which nearly all tasks are assigned to the fastest facility. We therefore generated problems on three facilities in which the processing time is drawn from a uniform distribution on  $[2, 20]$ ,  $[2, 25]$  and  $[2, 30]$  for the three facilities, 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.

## 10 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. The CP problems were solved with the `assignAlternatives` and `setTimes` options, and the CP subproblems of the Benders method were solved with the `setTimes` option. Both result in substantially better performance.

Table 1 displays computational results for minimum cost and minimum makespan problems on 2, 3 and 4 facilities. The CP solver is consistently faster than MILP on minimum cost problems. However, CP is unable to solve most problems with more than 16 tasks within two hours of computation time.

The Benders method is substantially faster than both CP and MILP. Its advantage increases rapidly with problem size, reaching some three orders of magnitude relative to CP for 16 or 18 tasks. Presumably the advantage would be greater for larger problems.

On minimum makespan problems, CP and MILP appear to be roughly competitive overall, with CP generally faster on 2-facility problems and MILP on 3 and 4-facility problems. Again the Benders method is considerably faster than both CP and MILP, with speedups ranging between two and three orders of magnitude for problems small enough to be solved by CP or MILP. As before, the speedup increases with the problem size.

As the number of tasks increases, the Benders subproblems in minimum cost and makespan problems eventually reach a size at which the computation time for the scheduling subproblem explodes. This point is reached later when there are more facilities, since the

subproblems are smaller when the tasks are spread over more facilities. In practice, precedence constraints or other side constraints could accelerate the solution of the scheduling subproblems and allow the solution of larger problems.

Table 2 investigates how the Benders method scales up to a larger number of facilities when the average number of tasks per facility is fixed to 5. The random instances are generated so that the fastest facility is roughly twice as fast as the slowest facility, with the other facility speeds spaced evenly in between. Since the subproblems remain relatively small as the problem size increases, it is possible to accommodate more tasks than in Table 1. In this case, the master problem (rather than the subproblem) dominates the computational task in larger problems and eventually becomes too hard to solve.

Table 2 suggests that the Benders method scales up better for minimum cost problems than for minimum makespan problems. Yet even when it fails to solve a makespan problem optimally, it obtains a feasible solution and a fairly tight lower bound on the optimal value. The bound steadily improve as the algorithm runs. Table 3 displays the bounds obtained for the minimum makespan problems of Table 2 that were not solved to optimality.

Table 4 shows computational results for minimum tardiness problems on three facilities using CP, MILP and Benders. Since problem difficulty tends to increase rapidly with minimum tardiness, the table shows results for individual problem instances and indicates the minimum tardiness for each. The instances are listed in order of minimum tardiness for each problem size.

The minimum tardiness problems are more readily solved by MILP than CP. Nonetheless the Benders method is faster than MILP for nearly all problem instances. The speedup is generally about one order of magnitude (actually the average ratio of solution times is 20.7 on problems that Benders solved optimally). The speedup is not nearly as great as for the minimum cost and minimum makespan problems, but significant nonetheless. These results were obtained by using the Benders cuts in (24), which rely on the solution of several modified subproblems in each iteration. As noted earlier, these are more effective than the cuts in (22), which resulted in performance that was almost always worse than MILP.

Despite its better performance, the Benders method does not significantly enlarge the class of tardiness problems that can be solved to optimality. On problems too large solve optimally, however, it finds much better solutions than MILP finds within a fixed amount of time (two hours). The Benders method also finds a lower bound on the optimal value that steadily increases with the number of iterations. The lower bounds are rather weak for the instances solved here, but they are tighter than the lower bound of zero obtained from the continuous relaxation of the MILP model.

## 11 Conclusions and Future Research

We find that logic-based Benders decomposition can improve substantially on the state of the art for the solution of planning and scheduling problems. In the case of minimum cost and minimum makespan problems, the Benders method is several orders of magnitude faster than either CP or MILP and solves larger problem instances to optimality.

On minimum tardiness problems, Benders decomposition tends to be about an order of magnitude faster than MILP (which is faster than CP on these problems). It does not solve significantly larger problems to optimality but obtains much better solutions than MILP for instances too large to solve optimally.

In fact, an important advantage of the Benders approach is that it can be terminated early while still yielding a feasible solution, as well as a lower bound on the optimal value that improves steadily as the algorithm runs. (This does not apply to the minimum cost problems, for which all intermediate solutions are infeasible.)

Future research should address several issues. One is whether effective Benders cuts can be developed for minimum makespan and minimum tardiness problems in which tasks have different release dates. In addition, implementation of a branch-and-check approach to solving the master problem could significantly improve performance, particularly on problems in which the master problem dominates the computational burden.

A fundamental issue is whether access to “dual” information from the CP scheduler (results of edge finding, etc.) would result in more effective Benders cuts, and how these cuts would be derived. This is particularly important for the tardiness problem, where the Benders cuts used here are rather weak.

A related issue is how to incorporate strong cuts in the master problem without generating a large number of constraints. Computational experience with the tardiness problems, for instance, suggests that stronger cuts could have been generated with little additional effort by solving several modified subproblems in each iteration. Each of the strong cuts, however, involves a logical condition that requires a large number of linear inequalities (many with “big-M” coefficients) in order to formulate it in the MILP master problem. Solution of the master problem bogs down as the cuts accumulate.

This suggests that rather than solving the master problem in the form of an MILP, it may be more effective to apply an algorithm that accommodates more complex logical constraints. This in combination with a branch-and-check strategy could result in a more robust class of Benders methods.

## 12 References

Benders, J. F. 1962. Partitioning procedures for solving mixed-variables programming problems, *Numerische Mathematik* 4: 238–252.

- Eremin, A., and M. Wallace. 2001. Hybrid Benders decomposition algorithms in constraint logic programming, in T. Walsh, ed., *Principles and Practice of Constraint Programming (CP 2001)*, *Lecture Notes in Computer Science* **2239**, Springer.
- Geoffrion, A. M. 1972. Generalized Benders decomposition, *Journal of Optimization Theory and Applications* **10**: 237–260.
- Hooker, J. N. 2000. *Logic-based Methods for Optimization: Combining Optimization and Constraint Satisfaction*, John Wiley & Sons.
- Hooker, J. N. and G. Ottosson. 2003. Logic-based Benders decomposition, *Mathematical Programming* **96**: 33–60.
- Hooker, J. N., and Hong Yan. 1995. 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) 267–288.
- Jain, V., and I. E. Grossmann. 2001. Algorithms for hybrid MILP/CP models for a class of optimization problems, *INFORMS Journal on Computing* **13**: 258–276.
- Thorsteinsson, E. S. 2001. Branch-and-Check: A hybrid framework integrating mixed integer programming and constraint logic programming, *Lecture Notes in Computer Science* **2239**: 16–30.
- Türkay, M., and I. E. Grossmann. 1996. Logic-based MINLP algorithms for the optimal synthesis of process networks, *Computers and Chemical Engineering* **20**: 959–978.



Table 1: Computation times in seconds for minimum cost and minimum makespan problems, using MILP, CP, and logic-based Benders methods. Each time represents the average of 5 instances. Computation was cut off after two hours (7200 seconds), and a + indicates that this occurred for at least one of the five problems.

| Facilities | Tasks | Min cost           |       |         | Min makespan |       |         |
|------------|-------|--------------------|-------|---------|--------------|-------|---------|
|            |       | MILP               | CP    | Benders | MILP         | CP    | Benders |
| 2          | 10    | 1.9                | 0.14  | 0.09    | 3.4          | 0.80  | 0.08    |
|            | 12    | 199                | 2.2   | 0.06    | 12           | 4.0   | 0.39    |
|            | 14    | 1441               | 79    | 0.04    | 2572+        | 299   | 7.8     |
|            | 16    | 3604+ <sup>a</sup> | 1511  | 1.1     | 5794+        | 3737  | 30      |
|            | 18    |                    | 7200+ | 7.0     |              | 7200+ | 461     |
|            | 20    |                    |       | 85      |              |       | 2656    |
|            | 22    |                    |       | 1674+   |              |       |         |
| 3          | 10    | 0.86               | 0.13  | 0.37    | 3.9          | 0.85  | 0.06    |
|            | 12    | 797                | 2.6   | 0.55    | 12           | 7.5   | 0.3     |
|            | 14    | 114                | 35    | 0.34    | 524          | 981   | 0.7     |
|            | 16    | 678 <sup>a</sup>   | 1929  | 4.5     | 1716+        | 4414  | 6.5     |
|            | 18    |                    | 7200+ | 15      | 4619+        | 7200+ | 13.3    |
|            | 20    |                    |       | 2.9     |              |       | 34      |
|            | 22    |                    |       | 23      |              |       | 3084+   |
| 4          | 10    | 2.0                | 0.10  | 0.6     | 1.0          | 0.07  | 0.09    |
|            | 12    | 7.2                | 1.4   | 4.0     | 5.0          | 1.9   | 0.09    |
|            | 14    | 158                | 72    | 2.8     | 24           | 524   | 0.8     |
|            | 16    | 906 <sup>a</sup>   | 344   | 0.8     | 35           | 3898  | 0.9     |
|            | 18    |                    | 6343+ | 5.2     | 3931+        | 7200+ | 13.9    |
|            | 20    |                    |       | 2.6     |              |       | 25      |
|            | 22    |                    |       | 22      |              |       | 472     |
| 24         |       |                    | 114   |         |              | 1131  |         |
| 26         |       |                    | 76    |         |              |       |         |

<sup>a</sup>CPLEX ran out of memory on one or more problems, which are omitted from the average time.

Table 2: Computation times in seconds for minimum cost and minimum makespan problems, using the Benders method. Each time represents the average of 5 instances. Computation was cut off after two hours (7200 seconds), and a + indicates that this occurred for at least one of the five problems.

| Tasks | Facilities | Min cost        | Min makespan |
|-------|------------|-----------------|--------------|
| 10    | 2          | 0.1             | 0.2          |
| 15    | 3          | 0.7             | 1.6          |
| 20    | 4          | 50 <sup>a</sup> | 13           |
| 25    | 5          | 2.9             | 213          |
| 30    | 6          | 4.8             | 3373+        |
| 35    | 7          | 128             | 6404+        |
| 40    | 8          | 1792+           | 7200+        |

<sup>a</sup>Includes one outlier that ran for 240 seconds.

Table 3: Best solution value and lower bound found after 7200 seconds of computation by the Benders method on the minimum makespan problems of Table 2 that were not solved to optimality.

| Tasks | Facilities | Best solution value | Lower bound |
|-------|------------|---------------------|-------------|
| 30    | 6          | 13                  | 12          |
| 35    | 7          | 11                  | 10          |
| 35    | 7          | 15                  | 13          |
| 40    | 8          | 14                  | 11          |
| 40    | 8          | 15                  | 12          |
| 40    | 8          | 16                  | 13          |
| 40    | 8          | 10                  | 9           |
| 40    | 8          | 13                  | 11          |

Table 4: Computational results for minimum tardiness problems on three facilities. Computation is terminated after two hours (7200 seconds).

| Tasks | CP     | Time (sec) |         | Best solution value found <sup>1</sup> |           | Benders lower bound <sup>2</sup> |
|-------|--------|------------|---------|--|-----------|----------------------------------|
|       |        | MILP       | Benders | MILP                                   | Benders   |                                  |
| 10    | 13     | 4.7        | 2.8     | <b>10</b>                              | <b>10</b> |                                  |
|       | 1.1    | 6.4        | 1.5     | <b>10</b>                              | <b>10</b> |                                  |
|       | 1.4    | 6.4        | 1.9     | <b>16</b>                              | <b>16</b> |                                  |
|       | 4.6    | 32         | 4.5     | <b>17</b>                              | <b>17</b> |                                  |
|       | 8.1    | 33         | 23      | <b>24</b>                              | <b>24</b> |                                  |
| 12    | 4.7    | 0.7        | 0.3     | <b>0</b>                               | <b>0</b>  |                                  |
|       | 14     | 0.6        | 0.3     | <b>0</b>                               | <b>0</b>  |                                  |
|       | 25     | 0.7        | 0.2     | <b>1</b>                               | <b>1</b>  |                                  |
|       | 19     | 15         | 2.3     | <b>9</b>                               | <b>9</b>  |                                  |
|       | 317    | 25         | 11      | <b>15</b>                              | <b>15</b> |                                  |
| 14    | 838    | 7.0        | 2.2     | <b>1</b>                               | <b>1</b>  |                                  |
|       | 7159   | 34         | 4.2     | <b>2</b>                               | <b>2</b>  |                                  |
|       | 1783   | 45         | 20      | <b>15</b>                              | <b>15</b> |                                  |
|       | > 7200 | 73         | 47      | <b>19</b>                              | <b>19</b> |                                  |
|       | > 7200 | > 7200     | 3123    | 26                                     | <b>26</b> |                                  |
| 16    | > 7200 | 19         | 1.2     | <b>0</b>                               | <b>0</b>  |                                  |
|       | > 7200 | 46         | 1.6     | <b>0</b>                               | <b>0</b>  |                                  |
|       | > 7200 | 52         | 5.3     | <b>4</b>                               | <b>4</b>  |                                  |
|       | > 7200 | 1105       | 187     | <b>20</b>                              | <b>20</b> |                                  |
|       | > 7200 | 3424       | 752     | <b>31</b>                              | <b>31</b> |                                  |
| 18    |        | 187        | 4.0     | <b>0</b>                               | <b>0</b>  |                                  |
|       |        | 15         | 8.1     | <b>3</b>                               | <b>3</b>  |                                  |
|       |        | 46         | 53      | <b>5</b>                               | <b>5</b>  |                                  |
|       |        | 256        | 54      | <b>11</b>                              | <b>11</b> |                                  |
|       |        | > 7200     | 1146    | 14                                     | <b>11</b> |                                  |
| 20    |        | 105        | 11      | <b>0</b>                               | <b>0</b>  |                                  |
|       |        | 4141       | 16      | <b>1</b>                               | <b>1</b>  |                                  |
|       |        | 39         | 28      | <b>4</b>                               | <b>4</b>  |                                  |
|       |        | 1442       | 305     | <b>8</b>                               | <b>8</b>  |                                  |
|       |        | > 7200     | > 7200  | 75                                     | 37        | 8                                |
| 22    |        | 6.3        | 20      | <b>0</b>                               | <b>0</b>  |                                  |
|       |        | 584        | 36      | <b>2</b>                               | <b>2</b>  |                                  |
|       |        | > 7200     | > 7200  | 120                                    | 40        | 7                                |
|       |        | > 7200     | > 7200  | 162                                    | 46        | 11                               |
|       |        | > 7200     | > 7200  | 375                                    | 128       | 34                               |
| 24    |        | 10         | 661     | <b>0</b>                               | <b>0</b>  |                                  |
|       |        | > 7200     | 53      | 20                                     | <b>0</b>  |                                  |
|       |        | > 7200     | 72      | 57                                     | <b>0</b>  |                                  |
|       |        | > 7200     | > 7200  | 20                                     | 5         | 3                                |
|       |        | > 7200     | > 7200  | 25                                     | 7         | 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.