A Constraint Programming Approach for a Batch Processing Problem with Non-identical Job Sizes

Arnaud Malapert^{a,c,*}, Christelle Guéret^b, Louis-Martin Rousseau^c

^aÉcole des Mines de Nantes, LINA UMR CNRS 6241, Nantes, France ^bÉcole des Mines de Nantes, IRCCyN UMR CNRS 6597, Nantes, France

^cÉcole Polytechnique de Montréal, CIRRELT, Montréal, Québec, Canada

Abstract

This paper presents a constraint programming approach for a batch processing machine on which a finite number of jobs of non-identical sizes must be scheduled. A parallel batch processing machine can process several jobs simultaneously and the objective is to minimize the maximal lateness. The constraint programming formulation proposed relies on the decomposition of the problem into finding an assignment of the jobs to the batches, and then minimizing the lateness of the batches on a single machine. This formulation is enhanced by a new optimization constraint which is based on a relaxed problem and applies cost-based domain filtering techniques. Experimental results demonstrate the efficiency of cost-based domain filtering techniques. Comparisons to other exact approaches clearly show the benefits of the proposed approach: it can optimally solve problems that are one order of magnitude greater than those solved by a mathematical formulation or by a branch-and-price.

Keywords: combinatorial optimization, artificial intelligence, constraint programming, scheduling, packing.

1. Introduction

This paper presents a constraint programming approach for a batch processing machine on which a finite number of jobs with non-identical sizes must be scheduled. A parallel batch processing machine can process several jobs simultaneously. Such machines are encountered in chemical, pharmaceutical, aeronautical, and semiconductor wafer industries where an oven, a drier, or an autoclave is used during the process. For example the composite components used in aerospace are made up of a matrix covered by carbon fiber which are consolidated into a solid structure at elevated temperature and pressure in an autoclave. Several parts of different size can be processed in a same autoclave at the same time. Another example is the final stress testing stage of semiconductor

christelle.gueret@mines-nantes.fr (Christelle Guéret), louis-martin.rousseau@polymtl.ca (Louis-Martin Rousseau)

April 5, 2012

 $^{^{*}}$ Corresponding author

Email addresses: arnaud.malapert@mines-nantes.fr (Arnaud Malapert),

Preprint submitted to Elsevier

manufacturing, where several different integrated circuits are transferred to an oven at elevated temperature for an extended period of time.

Early work on batch processor models can be traced to Ikura & Gimple (1986) who proposed an optimal algorithm for a problem with identical job processing times, identical job sizes, dynamic job arrivals and the objective of minimizing the makespan. Since then, heuristics (Perez et al., 2000; Wang & Uzsoy, 2002; Uzsoy, 1995), genetic algorithm (Wang & Uzsoy, 2002) and exact methods (Webster & Baker, 1995; Mehta & Uzsoy, 1998; Liu et al., 2007) have been proposed for identical job sizes and due date related performance measures. Several papers describe approaches for batch processing machine problems with non-identical job sizes, but concern mostly completion time related performance measures: heuristics (Azizoglu & Webster, 2000); genetic algorithm (Damodaran et al., 2006); simulated annealing (Damodaran et al., 2007); exact methods (Azizoglu & Webster, 2000, 2001; Dupont & Dhaenens-Flipo, 2002; Parsa et al., 2010; Kashan et al., 2009; Sabouni & Jolai, 2010). For an extensive review on scheduling with batching, we refer the reader to Potts & Kovalyov (2000).

On the other hand, constraint programming (CP) is an appealing technology in a variety of combinatorial problems which has grown steadily since the last three decades. According to Boucher et al. (1997), the main reason for success of CP on scheduling problems is the use of constraint propagation, which aims at removing from variable domains combinations of values which cannot appear in any consistent solution. Surprisingly, although CP has been successfully used to solve various scheduling problems, only a few papers concern scheduling problems with batching decisions. Moreover, these works focus on *serial*-batching problems in which a serial batching machine can process jobs contiguously as a batch whereas we discuss a *parallel*-batching problem. For example, Zeballos & Henning (2003); Kotecha et al. (2007); Felizari et al. (2009) discuss the modeling aspects of the multistage batch scheduling problem where each product has to be sequentially processed in a number of stages, and each stage contains one or more parallel processing units. The methods proposed in these papers assume that batching and scheduling decisions are made independently, *i.e.* each product order is divided into a number of batches (batching), which are then assigned to processing units and sequenced (scheduling). In another paper, Vilím (2007) proposes new filtering algorithms for a *serial*-batching machine with job families and sequence dependent setup times which adjust the time windows of the tasks.

In this paper, we propose a new constraint programming approach for a problem derived from a real application in the aeronautical industry. Composite components used in this kind of industry are made up of two categories of constituent materials: matrix (epoxy resin) and reinforcement (carbon fiber). They are fabricated according to the following process: the matrix is covered by carbon fiber, then the resulting part is consolidated into a solid structure at elevated temperature and pressure in an autoclave. Several parts can be processed at the same time as a batch in a same autoclave (batch processing machine). The size of a batch is then limited by the capacity (volume) of the autoclave. The processing time of each part (job) in the autoclave partially depends on the size of the part. When several parts are regrouped in an autoclave, the processing time of the batch is the longest processing time of its parts. As the processing times of such operations are longer than the other operations of the process, theses machines are often bottleneck. To our knowledge, only two papers (Daste et al., 2008b,a) concern the resolution of the problem which consists in scheduling jobs with non-identical job sizes on a batch processing machine to minimize the maximal lateness. Daste et al. propose a mathematical formulation and a branch-and-price.

More formally, the problem can be described as follows. A set J of n jobs and one single parallel batch processing machine with capacity b are given. Each job j is characterized by an integer triplet (p_j, s_j, d_j) , where p_j is its processing time, s_j is its size, and d_j is its due date. The sizes of the jobs are non-identical. The batch processing machine can process several jobs simultaneously as a batch as long as the sum of the sizes of the jobs that are in the batch does not exceed the capacity b of the machine. The processing time of a batch is equal to the longest processing time among the jobs in the batch. The completion time C_j of a job j is the completion time of the batch to which it belongs. The machine and jobs are assumed to be continuously available from time zero onward, or equivalently, they have equal release dates. Once the processing of a batch has been initiated, no job can be removed from or added to the batch. The objective is to minimize the maximal lateness $L_{max} = \max_{1 \le j \le n} (C_j - d_j)$. This problem, denoted by $1|p\text{-batch}; b < n; non\text{-identical}|L_{max}$, is unary NP-hard because Brucker et al. (1998) proved that the same problem with identical job sizes is unary NP-hard.

This paper presents the first constraint programming approach for the problem 1|p-batch; $b < n; non-identical|L_{max}$. After a brief overview of constraint programming techniques in Section 2, a new constraint programming model based on the decomposition of the problem is introduced in Section 3. Then, Section 4 describes a new optimization constraint based on the resolution of a relaxed problem enhanced by cost-based domain filtering techniques. In Section 5, specialized search strategies inspired from well- known bin packing approaches are presented. Finally, Section 6 evaluates the performance of filtering rules and search strategies, then compares our results to those of a mathematical formulation and to a branch-and-price approach.

2. Constraint programming background

Constraint programming (CP) has attracted high attention among experts from many areas because of its potential for solving hard real-life problems. A constraint satisfaction problem (CSP) consists of a set \mathcal{V} of variables defined by a corresponding set of possible values (the domains \mathcal{D}) and a set \mathcal{C} of constraints. A *constraint* is simply a logical relation between a subset of variables, each taking a value in its domain. The constraints thus restrict the possible values that variables can take. The important feature of constraints is their declarative manner, *i.e.* they only specify what relationship must hold. A partial assignment represents the case where the domains of some variables have been reduced, possibly to a singleton (namely a variable has been assigned a value). More formally, a partial assignment \mathcal{A} is the Cartesian product of all variable current domains. The current domain $\mathcal{D}(x)$ of each variable $x \in \mathcal{V}$ is always a (non-strict) subset of its initial domain. A partial assignment \mathcal{A}' extends another partial assignment \mathcal{A} , denoted $\mathcal{A}' \subseteq \mathcal{A}$, if the domain of any variable x in \mathcal{A}' is a subset of its domain in \mathcal{A} . Clearly, the relation \subseteq defines a partial order on partial assignments. A solution of a CSP is an assignment of a value to each variable such that all constraints are simultaneously satisfied (or consistent).

Solutions can be found by searching systematically through the possible assignments of values to variable. A *backtracking scheme* incrementally extends a partial assignment \mathcal{A} that specifies consistent values for some of the variable, toward a complete solution, by

repeatedly choosing a consistent value for another variable. The variables are labelled (given a value) sequentially and as soon as all the variables relevant to a constraint are assigned values, the validity of the constraint is checked. The variable and value selection heuristics select the next variable and value to try. If a partial solution violates any of the constraints, backtracking is performed to the most recently assigned variable that still has alternative values available in its domain. Clearly, whenever a partial assignment violates a constraint, backtracking is able to eliminate a subspace from the Cartesian product of all variable domains. Note that the running complexity for most nontrivial problems is exponential.

Another approaches, called *consistency techniques*, consist in removing inconsistent values from the domains of the variables until a solution is found. A filtering algorithm is associated to each constraint which removes inconsistent values from the domains of the variables, *i.e.* assignments which can not belong to a solution of the constraint. Constraints are handled through a *constraint propagation mechanism* which allows the reduction of the domains of variables until a global fixpoint is reached (no more domain reductions are possible). Different consistency levels are often available for a constraint. In fact, a constraint specifies what relationship must hold and its filtering algorithm is the computational procedure which enforces that relationship. Generally, consistency techniques are not complete, *i.e.* they do not remove all inconsistent values from the domains of the variables. Since each constraint applies its own local consistency techniques, global inconsistencies, *i.e.* associated to several constraints, are not always detected. Global constraints lessen this drawback by using the semantic information from subproblems in order to detect more inconsistencies or to reduce the computation time. Global constraints often exploits complex combinatorial structures such as graphs (cliques, matching, trees, paths \ldots).

Both systematic search and (some) consistency techniques can be used alone to solve a CSP completely, but their combination allows the search space to be explored in a complete and more efficient way. In this case, the propagation mechanism allows the reduction of the domains of variables and the pruning of the search tree. Scheduling is probably one of the most successful areas for CP thanks to specialized global constraints, which allow modelling an expressive and concise condition involving a non-fixed number of variables, as for instance resource limitations. For an extensive review on constraint programming and constraint-based scheduling, we refer the reader to Rossi et al. (2006) and Baptiste et al. (2001) respectively.

3. Constraint programming formulation

Our constraint programming formulation relies on the decomposition of the problem into finding an assignment of the jobs to the batches, and then minimizing the maximal lateness of the batches on a single machine. The problem of assigning the jobs to the batches is equivalent to the one-dimensional bin packing problem. Indeed, the definition of this problem is the following: given n indivisible items (jobs), each of a known nonnegative size s_j , and m bins (batches), each of capacity b, can we pack the n items into the m bins such that the sum of the sizes of the items in any bin is not greater than b? This problem has been shown NP-complete (Garey & Johnson, 1979). Then, once the jobs are packed into the batches, the problem of scheduling the batches is equivalent to minimizing the maximal lateness of a set of jobs (batches) on a single machine. This problem, denoted as $1||L_{max}$, is polynomially solvable (Lawler, 1973): an optimal schedule is obtained by applying Jackson's scheduling rule, also known as the earliest due date (EDD-)rule which schedules the tasks in order of non decreasing due dates. The main advantages of such a decomposition are twofold. First, the search space is restricted to batching decisions because sequencing decisions are independently computed in polynomial time. Then, the new optimization constraint introduced in Section 4 strongly relies on the decomposition because its filtering rules are based on a relaxation of the batch sequencing subproblem.

We now present the CP model of the studied problem. Without loss of generality, we assume in the remainder of the paper that the jobs are numbered according to non-increasing size $(s_j \geq s_{j+1})$, and that the size of each job is not greater than the batch capacity $(s_j \leq b)$. Note that the number m of batches is lower than or equal to n. Let J = [1, n] denote the set of job's indices and K = [1, m] denote the set of batch's indices. Let $d_{max} = \max_J \{d_j\}$ and $p_{max} = \max_J \{p_j\}$ be respectively the greatest due date and processing time of the jobs. Let $B_j \in K$ denote the batch where the job j is packed, and $J_k \subseteq J$ denote the set of jobs which are packed into the batch k. These variables satisfy the relation: $\forall j \in J$, $B_j = k \Leftrightarrow j \in J_k$. The non-negative integer variables $P_k \in [0, p_{max}]$, $D_k \in [0, d_{max}]$ and $S_k \in [0, b]$ represent the processing time, the due date, and the load of the batch k respectively. Lastly, let $M \in [0, m]$ and L_{max} be the number of non-empty batches and the objective variable (unbounded) respectively. The constraint programming formulation is given below:

 $\max \text{OfASet}\left(P_k, J_k, \left[p_j\right]_J, 0\right) \qquad \forall k \in K$ (1)

- $\min \text{OfASet}\left(D_k, \ J_k, \ [d_j]_J, \ d_{max}\right) \qquad \forall k \in K$ (2)
- $\operatorname{pack}\left(\left[J_k\right]_K, \left[B_j\right]_J, \left[S_k\right]_K, M, \left[s_j\right]_J\right)$ (3)

sequenceEDD
$$([B_j]_I, [D_k]_K, [P_k]_K, M, L_{max})$$
 (4)

Constraints (1), where $[p_j]_J$ is the associative array of processing times, enforce that the duration P_k of batch k is the maximal duration of its jobs (set J_k) if the batch is not empty, and equals 0 otherwise. Similarly, Constraints (2) enforce that the due date D_k of a batch k is the minimal due date of its jobs if the batch is not empty, and equals d_{max} otherwise. Indeed, the lateness of batch k defined as $\max_{j \in J_k} (C_k - d_j)$ is equal to $C_k - \min_{j \in J_k} (d_j)$ where C_k is the completion time of the batch k. Note that the lateness is negative if the job is early, nil if the job is on time, and positive if the job is tardy.

Constraint (3) is inspired from the global constraint of Shaw (2004) for the bin-packing problem. This constraint uses propagation rules incorporating knapsack-based reasoning, as well as a dynamic lower bound on the number of non-empty bins. This constraint replaces the channeling constraints between assignment variables $(\forall j \in J, B_j = k \Leftrightarrow j \in J_k)$ and enforces the consistency between assignments and loads $(\forall k \in K, \sum_{j \in J_k} s_j = S_k)$. Furthermore, pack propagates the redundant constraint specifying that the sum of the bin loads is equal to the sum of the item sizes $(\sum_{j \in J} s_j = \sum_{k \in K} S_k)$. Note that the limited capacity of the batches is enforced by the initial domain of the load variables S_k . AppendixA page 19 describes our implementation of pack.

Lastly, Constraint (4) enforces that the objective value L_{max} is equal to the maximal lateness of the batches scheduled according to the EDD-rule. This constraint applies several filtering rules that are explained in details in Section 4.

Finally, a solution to the constraint programming formulation is composed of a feasible

assignment of the jobs to the batches, and of the maximal lateness of the instance of the problem $1||L_{max}$ associated with these batches. Note that this model is also valid in the presence of additional constraints such as heterogeneous capacities (restrictions on the domains $\mathcal{D}(S_k)$), job incompatibilities (inequalities between variables B_j), load balancing (a spread constraint over variables S_k imposes a maximum standard deviation and an interval for the mean).

Dominance conditions on feasible assignments, which allow to consider a small subset of them, are a key property to solve bin packing problems (Martello & Toth, 1990; Shaw, 2004; Fukunaga & Korf, 2007). These conditions often consider that equal-sized items and equal-loaded bins can be swapped without sacrificing solution quality. These rules cannot be applied in batching machine problems because swapping equal-sized items or equal-loaded bins may modify the durations and due dates of the batches and possibly sacrificing solution quality. However, the search space still can be reduced by making sure that two jobs i and j such that $s_i + s_j > b$ belong to different batches. Since the jobs are sorted according to non-increasing size, let j_0 denote the largest index such that any pair of jobs with indices lower than j_0 are in different batches: $j_0 = \max\{j \mid \forall 1 \le i < j, s_{i+1} + s_i > b\}$. Then, Constraints (5) which pack the largest jobs into the first consecutive batches, can be added to the formulation.

$$B_j = j \qquad 1 \le j \le j_0 \tag{5}$$

4. Description of the sequenceEDD constraint

Pruning generally derives from feasibility reasoning. When coping with optimization problems, pruning can be done also on the basis of costs, *i.e.* optimality reasoning. Propagation can be aimed at removing combination of values which cannot lead to solutions whose cost is better than the best one found so far. Focacci et al. (1999) proposed to embed in global constraints optimization components representing suitable relaxations of the constraint itself. These components provide efficient operations research algorithms computing the optimal solution of the relaxed problem and a gradient function representing the estimated cost of each variable-value assignment. They show the benefit of using this information for pruning and for guiding the search on a variety of combinatorial optimization problems. In this section, we introduce an optimization constraint following the idea of Focacci et al. based on new relaxations and gradient functions.

As explained in Section 3, the global constraint sequenceEDD enforces that the objective value L_{max} is equal to the maximal lateness of an EDD-sequence of batches. This constraint uses a relaxation of the problem that yields a lower bound for the objective function to prune portions of the search space. The general idea is to infer primitive constraints on the basis of information on costs. We use optimization components within a global constraint representing a proper relaxation of the problem, which consists in minimizing the maximal lateness of batches on a single machine. The optimization components provide the optimal solution of the relaxed problem, its value and a gradient function computing the cost to be added to the optimal solution for some variable-value assignments. The optimal value of this solution improves the lower bound of the objective function and prunes portions of the search space for which their lower bound is bigger than the best solution found so far. Section 4.1 defines a relaxed problem in which the jobs that have not yet been assigned to a batch are ignored. Then, four filtering rules are presented in Sections 4.2, 4.3 and 4.4.

4.1. Relaxed problem

In this section, we describe the relaxed instance $I(\mathcal{A})$ of the problem $1||L_{max}$ built upon a partial assignment \mathcal{A} of the variables, as well as an algorithm to solve it. At each point in the resolution, the relaxed problem consists in scheduling the current batches (partially filled) without considering the remaining jobs (not yet assigned to a batch). Let recall that, at each point in the resolution, we have a partial assignment \mathcal{A} which we define as the set of current domains of all variables. The current domain $\mathcal{D}(x)$ of a variable x is always a (non-strict) subset of its initial domain. An assignment \mathcal{A}' extends a partial assignment \mathcal{A} , denoted $\mathcal{A}' \subseteq \mathcal{A}$, if the domain of any variable x in \mathcal{A}' is a subset of its domain in \mathcal{A} . Let $\min(x)$ and $\max(x)$ be the minimum and maximum value of the domain $\mathcal{D}(x)$ in the current partial assignment. Let $x \leftarrow v$ denote the restriction of the domain of x to a single value v. Let $\delta_1, \delta_2, \ldots, \delta_{n^*}$ be the distinct increasing values of the due dates d_j $(j \in J)$. Note that n^* can be smaller than n if some jobs have identical due dates. However, for sake of clarity, we will consider that $n^{\star} = n$. Let $K_{\mathcal{R}_q}(\mathcal{A}) = \{k \in K \mid \max(D_k) \mathcal{R} \ \delta_q\}$ be the set of batches related to the due date δ_q by the arithmetic relation $\mathcal{R} \in \{<,\leq,=,\geq,>\}$. Similarly, let $J_{\mathcal{R}q}(\mathcal{A}) = \{j \in J \mid d_j \ \mathcal{R} \ \delta_q\}$ be the set of jobs related to the due date δ_q . Finally, let $P(\mathcal{A}, \tilde{K}) = \sum_{k \in \tilde{K}} \min(P_k)$ be the minimal total duration of a set of batches $\tilde{K} \subseteq K$ in the partial assignment \mathcal{A} .

An instance $I(\mathcal{A})$ of the relaxed problem consists of n buckets where a bucket $q \in J$ is the set of batches $K_{=q}(\mathcal{A})$. Each bucket q has a due date δ_q and a processing time $\pi_q(\mathcal{A}) = P(\mathcal{A}, K_{=q}(\mathcal{A}))$. As the buckets are, by definition, numbered according to a strictly increasing order of their due date $(\delta_q < \delta_{q+1})$, minimizing the maximal lateness of $I(\mathcal{A})$ is equivalent to computing the maximal lateness of the sequence of buckets in this order. Let $C_q(\mathcal{A}) = \sum_{i=1}^q \pi_i(\mathcal{A})$ and $L_q(\mathcal{A}) = C_q(\mathcal{A}) - \delta_q$ be respectively the completion time and lateness of bucket q in the sequence. Therefore, the optimal objective value of instance $I(\mathcal{A})$ is given by: $L(\mathcal{A}) = \max_{q \in J} (L_q(\mathcal{A}))$.

Figure 1 illustrates the solutions of two relaxed problems for a batching machine of capacity b = 10. Table 1(a) contains the jobs to schedule in the original problem 1|p-batch; b < n; non-identical $|L_{max}$. Table 1(b) gives the buckets of instance $I(\mathcal{A}_1)$ where \mathcal{A}_1 is a solution of the problem, *i.e.* a total assignment. In this example, note that there are less buckets $(n^* = 3)$ than jobs (n = 4). Figure 1(d) shows an optimal solution of the relaxed instance $I(\mathcal{A}_1)$ which is also an optimal solution of the original problem. The batching machine is represented as a drawing where the horizontal and vertical axes correspond respectively to the time and the load of the resource. A job is represented as a rectangle for which the length and the height respectively match its duration and its size. Jobs with identical starting times belong to the same batch. The solution contains three batches: the first one contains only job 1; the second one contains the jobs 2 and 4, and has a load equal to $s_2 + s_4 = 9$, a minimal duration equal to $\max\{p_2, p_4\} = 9$ and a maximal due date equal to $\min\{d_2, d_4\} = 2$; the third one contains only job 3. A bucket is represented as a dashed rectangle which encapsulates its batches. The batches of a bucket are separated by a dashed line with a square marker. The lateness of each bucket is represented on the right part of the figure as a line starting from its due date and ending at its completion time. At this point, the bucket 1 contains batches 1 and 2



Figure 1: Two illustrative examples of the construction of $I(\mathcal{A})$.

because $\max(D_1) = \max(D_2) = d_1$. The bucket 2 is empty because there is no batch k such that $\max(D_k) = \delta_2$ and its lateness, drawn as a dashed line, is therefore dominated by its first non-empty predecessor (bucket 1). The bucket 3 contains batch 3 with job 3. Table 1(c) gives the buckets of instance $I(\mathcal{A}_2)$ where \mathcal{A}_2 is a partial assignment in which the job 4 is not yet assigned. Figure 1(e) shows a solution of the relaxed instance $I(\mathcal{A}_2)$ in which the job 4 does not appear. At this point, each bucket contains a unique batch: the first, second and third buckets contain respectively the batches 1, 2 and 3. In fact, \mathcal{A}_2 is inferred at the root node by the propagation of Constraints (5) which enforces that $B_1 = 1, B_2 = 2$, and $B_3 = 3$.

At each point of the resolution, the construction of instance $I(\mathcal{A})$ can be done in O(n). Indeed, the due dates of the jobs need to be sorted in increasing order which takes $O(n \log n)$, but this can be done once and for all before starting the search process. Then, building instance $I(\mathcal{A})$ consists in assigning each batch k to a bucket q. This can be done in O(1) as the bucket of a batch k is the bucket q such that $\delta_q = \max(D_k)$. Inserting batch k in its bucket, and updating the duration of the bucket is also done in constant time. The overall complexity is O(n), because the number of batches m is lower than n. Finally, the resolution of $I(\mathcal{A})$ simply consists in sequencing the buckets in their

numbering order, which can be done in O(n). The following sections present filtering rules based on the resolution of this relaxed problem.

4.2. Filtering the lateness variable

The filtering rules on the lateness variable are based on the two following propositions.

Proposition 1. L is a monotonic function from the partially ordered set of partial assignments onto the integers:

$$\mathcal{A}' \subseteq \mathcal{A} \Rightarrow L(\mathcal{A}') \ge L(\mathcal{A}).$$

Proof. Since Constraints (1) enforce that the minimal duration of a batch in \mathcal{A}' is greater than or equal to its minimal duration in \mathcal{A} , the total duration of a set of batches \tilde{K} does not decrease from \mathcal{A} to \mathcal{A}' : $\forall \tilde{K} \subseteq K, \ P(\mathcal{A}, \tilde{K}) \leq P(\mathcal{A}', \tilde{K}).$

Furthermore, since Constraints (2) enforce that assigning a new job in a batch belonging to bucket q can only imply its transfer to a bucket q' such that $q' \leq q$, the set of batches before bucket q can only increase from \mathcal{A} to \mathcal{A}' : $\forall q \in J$, $K_{\leq q}(\mathcal{A}) \subseteq K_{\leq q}(\mathcal{A}')$. Besides, since P_k is a non-negative integer variable, the function $P(\mathcal{A}, \tilde{K})$ is an increasing function from the sets of batches onto the integers: $\tilde{K} \subseteq \tilde{K}' \Rightarrow P(\mathcal{A}, \tilde{K}) \leq P(\mathcal{A}, \tilde{K}')$. Therefore, the following inequalities hold:

$$C_q(\mathcal{A}) = \sum_{1 \le i \le q} \pi_i(\mathcal{A}) = P(\mathcal{A}, K_{\le q}(\mathcal{A}))$$

$$\le P(\mathcal{A}', K_{\le q}(\mathcal{A})) \le P(\mathcal{A}', K_{\le q}(\mathcal{A}'))$$

$$< C_q(\mathcal{A}').$$

Since the completion time of each bucket q can only increase from \mathcal{A} to \mathcal{A}' , the monotonicity of the function $I(\mathcal{A})$ is proven as follows: $\forall q \in J, \ C_q(\mathcal{A}') \geq C_q(\mathcal{A}) \Rightarrow L(\mathcal{A}) \leq L(\mathcal{A}')$.

Proposition 2. Once all durations and due dates of the batches are fixed, the optimal schedule of the relaxed problem $I(\mathcal{A})$ corresponds also to an optimal schedule of any feasible assignment which extends \mathcal{A} .

Proof. Once all durations and due dates of the batches are assigned values, the durations and latenesses of buckets stay unchanged until a feasible solution is found. Therefore, the proposition is satisfied if and only if the maximal lateness of the relaxed problem is equal to the maximal lateness of any solution which extends \mathcal{A} . As all the batches of a bucket q have the same due date δ_q , all the orderings of these batches are equivalent regarding the EDD-rule. Thus, the optimal schedule of the buckets for the relaxed problem corresponds to an optimal schedule of the batches (respecting the EDD-rule) for the initial problem. Note that scheduling the batches of a bucket according to non-decreasing processing time improves the average lateness of these batches.

Let recall that $x \leftarrow v$ denotes the restriction of the domain of x to a single value v, let $\min(x) \leftarrow v$ denote the restriction to values greater than v, and let $\max(x) \leftarrow v$ denote the restriction to values lower than v. Furthermore, let $x \not\leftarrow v$ denote the elimination of a single value v. From Proposition 2, we deduce the *final lateness filtering rule* (FF)

which solves the relaxation problem to assign the objective variable once all durations and due dates have been assigned values:

$$\forall k \in K, P_k \text{ and } D_k \text{ are assigned} \Rightarrow L_{max} \leftarrow L(\mathcal{A}).$$
 (FF)

Note that this case occurs sometimes before a total assignment of the jobs to the batches is reached, since propagation can reduce all domains to singleton. A corollary of Propositions 1 and 2 is that, at each point of the search, the maximal lateness $L(\mathcal{A})$ of $I(\mathcal{A})$ is a lower bound on any feasible schedule which extends \mathcal{A} . Thus, the minimal objective value can be updated at each relevant domain change by the *lateness filtering rule* (LF):

$$\exists k \in K, \min(P_k) \text{ or } \max(D_k) \text{ have changed} \Rightarrow \min(L_{max}) \leftarrow L(\mathcal{A})$$
 (LF)

4.3. Cost-based domain filtering of assignments

The cost-based domain filtering rule of assignments (AF) reduces the search space based on the marginal cost of the assignment of job j to batch k ($B_j \leftarrow k$). The idea is to eliminate, at each relevant domain change, every job j, as a candidate for packing in a batch k, if the marginal cost associated with the assignment of the job j to the batch k exceeds the best upper bound found so far, or more formally:

$$\exists k \in K, \min(P_k) \text{ or } \max(D_k) \text{ have changed} \Rightarrow \\ \forall j \in J, \text{ such that } |\mathcal{D}(B_j)| > 1 \text{ and } \forall k \in \mathcal{D}(B_j), \\ L(\mathcal{A} \cap \{B_j \leftarrow k\}) > \max(L_{max}) \Rightarrow B_j \not\leftarrow k \quad (AF)$$

where $\mathcal{A} \cap \{B_j \leftarrow k\}$ stands for $\mathcal{A} \cap \{B_j \leftarrow k, \min(P_k) \leftarrow p_j, \max(D_k) \leftarrow d_j\}$. Indeed, the propagation of Constraints (1) and (2) after the assignment $B_j \leftarrow k$ implies respectively that $\min(P_k) \leftarrow p_j$ and $\max(D_k) \leftarrow d_j$.

A simple filtering algorithm can compute the marginal cost from scratch for each possible assignment with an overall complexity of $O(n^3)$. We propose, in AppendixB, an O(nm) version of this algorithm based on the incremental computation of marginal costs.

4.4. Cost-based domain filtering based on bin packing

In this section, we introduce a cost-based domain filtering rule based on the computation of marginal cost associated with a number of non-empty batches M. Let $K^* = \{k \in K \mid \exists j \in J, B_j = k\}$ be the set of non-empty batches. Note that we assume that $|K^*| = \max\{k \in K^*\}, i.e.$ jobs are packed into the first consecutive batches. Let $J^* = \{j \in J \mid \mathcal{D}(B_j) > 1 \land \max(B_j) > |K^*|\}$ denote the set of *open* jobs, *i.e.* unpacked jobs that can be used to create new batches. The rule (PF1) updates the possible number of non-empty batches by considering the current number of non-empty batches and the number of open jobs.

$$\min(M) \leftarrow |K^*| \qquad \max(M) \leftarrow |K^*| + |J^*| \tag{PF1}$$

This rule (also applied by pack) is required to ensure correctness of the reasoning presented below. Let $\mathcal{A}' = \mathcal{A} \cap (\bigcup_{j \in \tilde{J}} \{B_j \leftarrow k_j\})$ denote the assignment of a subset $\tilde{J} \subseteq J^*$ of open jobs to new batches, *i.e.* pairwise different empty batches $(\forall j \in \tilde{J}, k_j > |K^*|$ and $\forall i \neq j \in \tilde{J}, k_i \neq k_j)$. Therefore, the completion time and lateness of each bucket q have to be updated according to the total duration increase of its predecessors. Indeed, the new completion time of bucket q is equal to its completion time before the assignments plus the duration increase of its predecessors: $C_q(\mathcal{A}') = C_q(\mathcal{A}) + \sum_{j \in \tilde{J}_{\leq q}} p_j$. Unfortunately, the combinatorial of such assignments grows exponentially, and filtering them is difficult and costly.

Therefore, we compute a lower bound of $C_q(\mathcal{A}')$ by computing a lower bound of the completion time $C_q(\mathcal{A} \cap \{M \leftarrow m^*\})$ of a bucket q for a given number m^* of non-empty batches as follow. Let $\Pi(m^*)$ be the sum of m^* open jobs of smallest processing times. If $m^* - |K^*|$ new batches are created with the open jobs, at least $m^* - |K^*| - |J_{>q}^*|$ new batches are scheduled before bucket q. Then, the completion time of any bucket q after the creation of $m^* - |K^*|$ batches is greater than its completion time before the creation of the new batches plus the minimal sum $\Pi(m^* - |K^*| - |J_{>q}^*|)$ of the processing times of $m^* - |K^*| - |J_{>q}^*|$ open jobs:

$$C_q(\mathcal{A} \cap \{M \leftarrow m^*\}) = C_q(\mathcal{A}) + \Pi(m^* - |K^*| - |J_{>q}^*|)$$

$$\leq C_q\left(\mathcal{A} \cap \left(\bigcup_{j \in \tilde{J}}\{B_j \leftarrow k_j\}\right)\right) \qquad \forall \tilde{J} \subseteq J^*, \ |\tilde{J}| = m^* - |K^*|$$

As a consequence, the maximal lateness of any solution extending \mathcal{A} with exactly m^* nonempty batches is greater than $L(\mathcal{A} \cap \{M \leftarrow m^*\})$. The rule (PF2) updates the minimal objective value according to the marginal cost associated with the current minimum number of batches. The rule (PF3) decrements the maximum number of batches if its marginal cost exceeds the best upper bound found so far.

$$\min(L_{max}) \leftarrow L(\mathcal{A} \cap \{M \leftarrow \min(M)\}) \tag{PF2}$$

$$L(\mathcal{A} \cap \{M \leftarrow \max(M)\}) > \max(L_{max}) \Rightarrow \max(M) \leftarrow \max(M) - 1$$
 (PF3)

At each relevant domain change, the *cost-based domain filtering rule based on bin packing* (PF) applies the three filtering rules described above.

$$(\exists j \in J, \mathcal{D}(B_j) \text{ has changed}) \lor (\exists k \in K, \min(P_k) \text{ or } \max(D_k) \text{ have changed}) \Rightarrow$$

Apply rules (PF1), (PF2) and (PF3) (PF)

The computation of function Π is $O(n \log n)$ since it can be performed during initialization by sorting, and then summing the processing times of open jobs. The rules (PF1) and (PF2) are applied within linear time, since $L(\mathcal{A})$ can be computed in O(n). The rule (PF3) is applied until the marginal cost associated with the maximum number of batches becomes feasible, *i.e.* at most $|J^*|$ times. Therefore, the overall complexity of the cost-based domain filtering based on bin packing is $O(n^2)$.

Figure 2 illustrates the computation of marginal costs associated to the presence of exactly $m^* = 4$ non-empty batches for the instance introduced in Table 1(a) and an empty assignment \mathcal{A}_3 (before the initial propagation). The set of open jobs J^* is equal to the set of jobs J whereas the set of non-empty batches K^* is empty. Therefore, rules (PF1) and (PF2) do not modify the domain of M. Since the completion time $C_q(\mathcal{A}_3)$ of each bucket q is equal to 0, the new completion time $C_q(\mathcal{A}_3 \cap \{M \leftarrow 4\})$ is equal to $\Pi(m \star - |J_{>q}^*|)$. As a consequence, the bucket 1 ends at $\Pi(4 - |\{2,3\}|) = p_1 + p_3 = 12$, the bucket 2 ends at $\Pi(4 - |\{3\}|) = p_1 + p_3 + p_2 = 20$, and the bucket 3 ends at $\Pi(4 - |\emptyset|) = p_1 + p_3 + p_2 + p_4 = 29$.



Figure 2: Example of the computation of $L(\mathcal{A}_3 \cap \{M \leftarrow 4\})$ for the (empty) assignment \mathcal{A}_3 . $\mathcal{A}_3 = \{B_1 \in [1,4], B_2 \in [1,4], B_3 \in [1,4], B_4 \in [1,4]\}.$

If we suppose that the best upper bound found so far is equal to 15, then the maximum number of non-empty batches is reduced from 4 to 3 and the rule (PF2) is applied for $m^{\star} = 3$.

5. Search strategy

At each node in the search tree, the branching selects the variable B_j of a job j using a variable selection heuristic, and assigns job j to a batch k chosen with a value selection heuristic, *i.e.* takes the decision $B_j \leftarrow k$. On backtracking, the search states that the chosen job cannot be placed in the selected batch. If this batch becomes empty, a symmetry breaking rule is applied which eliminates "equivalent" batches, *i.e.* other empty batches, from the list of candidates for the current job. Note that this dynamic symmetry breaking rule generalizes the idea behind Constraints (5) because, on backtracking, we would forbid packing a large job into another batch.

Several variable selection heuristics based on durations, sizes and due dates can be used. We chose the heuristic called *complete decreasing* (Gent & Walsh, 1997) which packs jobs in order of non-increasing size. Indeed, preliminary experiments showed that packing large jobs first improves filtering of pack and sequenceEDD constraints more than simple variants using also processing times and due dates.

Concerning the value selection heuristics, we investigated two classical heuristics in bin packing problems, namely *first fit* that selects the first available batch, and *best fit* which selects an available batch with the least free space. We also propose a new value selection heuristic named *batch fit* which selects an available batch having the least value of $\gamma(B_j \leftarrow k)$ where $\gamma(B_j \leftarrow k)$ measures roughly the fitness of the assignment within the instance $I(\mathcal{A})$ as follows:

$$\gamma(B_j \leftarrow k) = \frac{|\min(P_k) - p_j|}{\max_J(p_j) - \min_J(p_j)} + \frac{|\max(D_k) - d_j|}{\max_J(d_j) - \min_J(d_j)}$$

The idea of this heuristic is that a perfect solution would be composed of batches containing jobs having all identical durations and due dates.

6. Experimental results

This section presents computational experiments conducted to evaluate our approach. Section 6.1 describes the set of instances on which the experiments have been done. Section 6.2 evaluates the performance of the different filtering rules. The performance of the

new value selection heuristic is studied in Section 6.3. Finally, our approach is compared to a mathematical formulation in Section 6.4 and to a branch-and-price in Section 6.5. All the experiments were conducted on a cluster of Linux machines, each node with 48GB of RAM and two quad core 2.4GHz processor. Our implementation is based on choco (http://choco.mines-nantes.fr) which is an open source java library for constraint programming built on an event-based propagation mechanism with backtrackable structures. The constraint pack has been integrated in the latest releases and the solver has been extended with the optimization constraint sequenceEDDand search heuristics. The constraint sequenceEDD always applies the algorithm for the rule (AF) presented in AppendixB. The time limit has been fixed to 3600 seconds (1h).

6.1. Instances

Our algorithm has been tested on randomly generated instances proposed by Daste et al. (2008b) ranging from 10 to 100 jobs. The processing times are determined using the uniform distribution $(p_j = U[1, 99])$. Job sizes are generated from discrete uniform distribution between 1 and 10, and the machine capacity is b = 10 (inspired by Ghazvini & Dupont, 1998). For a given instance, once the processing times and sizes of all jobs are computed, the due dates are generated using the following formula: $d_j = U[0, \alpha] \times \tilde{C}_{max} + U[1, \beta] \times p_j$ (inspired by Malve & Uzsoy, 2007) where $\tilde{C}_{max} = (\sum_{j=1}^n s_j \times \sum_{j=1}^n p_j) \div (b \times n)$ is an approximation of the time required to process all the jobs on the batch processing machine. For each number of jobs $n \in \{10, 20, 50, 75, 100\}$, 40 instances were generated with $\alpha = 0.1$ and $\beta = 3$.

6.2. Performance of the filtering rules

Improving the filtering of a constraint could benefit the resolution but it happens that simple and short algorithms outperform more complicated ones. Therefore, it is reasonable to evaluate their filtering power (see Section 6.2.1) and to ask whether cost-based domain filtering techniques are useful during search (see Section 6.2.2). In this section, the name of a filtering rule is an abbreviation which reflects its nested structure, *i.e.* (FF) \subset (LF) \subset (AF) \subset (PF). For instance, (PF) refers to the activation of the four filtering rules.

6.2.1. Destructive lower bound

We measured the performance of the different filtering rules by computing destructive lower bounds. Destructive lower bounds are obtained by first imposing an upper limit on the objective function value (say F) as low as possible (the reduced problem). Then, the technique consists in trying to contradict (destruct) quickly (using only propagation, and no search) the existence of a solution of value F. In case of success, then F is a valid lower bound, otherwise F is incremented by one unit and the contraction attempted again until the infeasibility of the a solution of value F can not be proven without searching. The computation of such a lower bound has the advantage of being independent from the search. Three different destructive lower bounds were computed using (LF), (AF), and (PF) respectively. Note that we ignore (FF) since it is unlikely that it performs any filtering in this case. Because destructive lower bounds are computed very quickly, another three (better) lower bounds were computed using also shaving techniques (see Rossi et al., 2006). Shaving is similar to prove by contradiction. We propagate the

	Init	Initial Propag.			Destr. LB			Destr. $LB + Shaving$					
n	\mathbf{LF}	AF	PF	\mathbf{LF}	AF	\mathbf{PF}	\mathbf{LF}	AF	\mathbf{PF}	$\bar{t}_{\rm LF}$	\bar{t}_{AF}	$\bar{t}_{\rm PF}$	
10 20 50	89.4 90.1 89.7	$89.4 \\ 90.1 \\ 89.7$	$89.7 \\ 90.4 \\ 90.2$	89.4 90.1 89.7	$97.9 \\ 95.6 \\ 93.6$	$98.1 \\ 95.7 \\ 93.6$	93.4 93.6 91.8	99.0 96.8 94.1	$99.2 \\ 97.0 \\ 94.1$	$\begin{array}{c c} 0.07 \\ 0.17 \\ 2.71 \end{array}$	0.05 0.15 2.71	$0.05 \\ 0.1 \\ 1.43$	
$\begin{array}{c} 75\\ 100 \end{array}$	88.5 87.2	$88.5 \\ 87.2$	$89.1 \\ 87.6$	88.5 87.2	$91.4 \\ 89.4$	$91.5 \\ 89.4$	89.8 88.3	$91.7 \\ 89.6$	$91.8 \\ 89.7$	8.8 23.20	$\begin{array}{c} 12.01 \\ 29.48 \end{array}$	$\begin{array}{c} 5.25\\ 14.04 \end{array}$	

Table 1: Quality of destructive lower bounds.

assignment of a job j to a batch k in the reduced problem. If an infeasibility is found, then the assignment is invalid and so we can eliminate job j as candidate for packing in batch k. In this case, if the propagation of the domain reduction $B_j \not\leftarrow k$ lead to a contradiction, then the reduced problem is infeasible. To limit computation time, shaving was used only once for each assignment $B_i \leftarrow k$.

Let ub denote the best upper bound found for a given instance reported in AppendixC. The quality of the computed lower bound lb for a given instance is evaluated by $(100 \times (lb + d_{max})) \div (ub + d_{max})$ (inspired by Malve & Uzsoy, 2007). This measure is equal to the optimality gap until the problem's size n = 20, since these instances have been solved optimally. Besides, it is almost equal to the optimality gap for problem's size n = 50, since only two instances are not solved optimally. Table 1 gives the average quality of the initial propagation (columns 2–4) and destructive lower bounds without and with shaving (columns 5–7 and 8–10) as a function of the number of jobs n (column 1). Average computation times \bar{t} (in seconds) are only reported for destructive lower bounds with shaving (columns 11–13) because they are not significant otherwise (\ll 1 second).

The rule (AF) does not strengthen the lower bound during the initial propagation, whereas it greatly improves the two destructive lower bounds. The rule (PF) slightly improves on (AF) in any case. Besides, destructive lower bounds without shaving are efficient and their computation times stay negligible. However, using shaving slightly improves the quality of destructive lower bounds, but leads to a significant increase of the computation time. The rule (PF) improves the destructive lower bound with shaving in comparison to the rule (AF) while reducing its computation time by half. Note that Constraints (5) greatly contribute to the quality of these lower bounds.

6.2.2. Complete decreasing first fit

The procedure *complete decreasing first fit* (Gent & Walsh, 1997) packs jobs in order of non-increasing size, packing each job in the first available batch that will accommodate it. Since the destructive lower bound mechanism is deactivated, the impact of the filtering on the decision process is lessened. We chose instances that could be solved optimally by most of the filtering rules, from (FF) to (PF), so that comparisons could be made in reasonable computation times. In this regard, we examined all instances with less than 50 jobs (120 instances).

Figure 3(a) shows the percentage of instances solved optimally as a function of the time in seconds. First, the performance of the rule (FF) shows that the model is able to capture some solutions even when filtering only happens once all durations and due dates have



Figure 3: Comparison of filtering rules using complete decreasing first fit.

been assigned values. We believe that the constraint programming approach, within this model, is extremely effective at identifying assignments that lead to equivalent EDD-schedule. Secondly, the rule (LF) based on the relaxed problem improves both solution quality and time over the final lateness rule (FF) but a number of instances remains unsolved optimally. Finally, only two instances with 50 jobs remain unsolved optimally within the time limit when using rule (AF) or (PF). However, this graph does not show clearly the gain offered by the rule (PF) over the rule (AF).

To this end, Figure 3(b) compares the resolution of rules (AF) and (PF). Each point represents one instance and its x coordinate is the ratio of the solving time with (AF) over the solving time with (PF), whereas its y coordinate is the ratio of the number of backtracks with (AF) over the number of backtracks with (PF). We only report the results of instances with solving times greater than 2 seconds, which happens only for problem of size n = 50. All points are located above and on the right of the point (1,1), because all instances are improved by the use of (PF). In fact, the rule (PF) lead to a similar behaviour than during the computation of destructive lower bounds. It always reduces the computation time and some instances are roughly identical using (AF) or (PF).

6.3. Performance of value selection heuristics

We evaluate the impact of value selection heuristics on the resolution. In this regard, we examine all instances using all filtering rules without applying any destructive lower bound. We only compare *batch fit* to *first fit*, because *first fit* and *best fit* gives equivalent results within our approach. Indeed, the difference between the solution times of *first fit* and *best fit* is always less than 5 seconds on instances with less than 50 jobs, and even in this case, the difference represents at most 2 % of the solution time. Furthermore, both heuristics give the same objective values for all instances with more than 50 jobs. Figure 4(a) analyses the effect of the value selection heuristic on instances with 50 jobs

Figure 4(a) analyses the effect of the value selection heuristic on instances with 50 jobs and solving times greater than 2 seconds. Each point represents one instance and its xcoordinate is the ratio of the solving time using *first fit* over the solving time using *batch*



Figure 4: Comparison of value selection heuristics.

fit, whereas its y coordinate is the ratio of the number of backtracks with first fit over the number of backtracks with batch fit. Note also that the scale is logarithmic, and that all points are around the diagonal because the number of backtracks is roughly proportional to the time (top-left and bottom-right quadrants are empty). All points located above or on the right of the point (1,1) are instances improved by the use of batch fit (top-right quadrant). The heuristic batch fit globally improves the solution time and even by a factor 4 on some instances. However, the performances over a very few instances located below (1,1) are slightly degraded (bottom-left quadrant).

Figure 4(b) analyses the effect of the value selection heuristic on instances with strictly more than 50 jobs. Let *lb* denote the best lower bound found for a given instance reported in AppendixC. The quality gap of a solution *ub* for a given instance is evaluated by $(ub + d_{max}) \div (lb + d_{max})$ (inspired by Malve & Uzsoy, 2007). This quality measure is often greater than the optimality gap since these instances have been infrequently solved optimally. Each point represents one instance and its *x* coordinate is the quality gap with *batch fit*, whereas its *y* coordinate is the quality gap with *first fit*. All points located above the line (x = y) are upper bounds improved by *batch fit*. Using *batch fit* globally improves the solution, but the performances over a few instances are slightly degraded.

6.4. Comparison to a mathematical formulation

Our approach is now compared to a mathematical formulation of the studied problem inspired by Daste et al. (2008b).

6.4.1. Mathematical formulation

Let x_{jk} be a boolean variable equal to 1 if the job j is assigned to the k-th batch of the sequence. The non-negative integer variables P_k , D_k , and C_k represent the processing time, due date, and completion time of the k-th batch respectively.

$$\min L_{max}$$
Subject to:
$$\sum_{k \in K} x_{jk} = 1 \qquad \forall j \in J \qquad (6)$$

$$\sum_{j \in J} s_j x_{jk} \leq b \qquad \forall k \in K \tag{7}$$

$$p_j x_{jk} \leq P_k \quad \forall j \in J, \ \forall k \in K$$
 (8)

$$C_{k-1} + P_k = C_k \qquad \forall j \in J, \ \forall k \in K \qquad (9)$$

$$(d_{max} - d_j)(1 - x_{jk}) + d_j \geq D_k \qquad \forall j \in J, \ \forall k \in K \qquad (10)$$

$$D_{k-1} \leq D_k \qquad \forall k \in K \tag{11}$$

$$C_k - D_k \leq L_{max} \quad \forall k \in K \tag{12}$$

$$\forall j \in J, \forall k \in K, x_{jk} \in \{0, 1\} \\ \forall k \in K, C_k > 0, P_k > 0, D_k > 0$$

The objective function aims to minimize the maximal lateness. Constraints (6) state that each job must be assigned to exactly one batch, and Constraints (7) ensure that no batch exceeds the machine capacity b. These constraints define the bin packing model proposed by Martello & Toth (1990). Constraints (8) state that the duration of each batch k is equal to the maximum duration of its jobs. Constraints (9) ensure that the completion time of the k-th batch is equal to the completion time of the (k-1)-th batch plus the processing time of the k-th batch, *i.e.* the solution is a sequence of batches in their numbering order without idle time. Note that the addition of a first fictitious batch ending at the initial time $(C_0 = 0)$ is required. Constraints (10) state that the due date D_k of the k-th batch is equal to the earliest due date of the jobs that are in the batch. Constraints (11) ensure that the sequence of batches satisfies the EDD-rule. These constraints help the resolution since most permutations of batches which lead to the same EDD-schedule are forbidden but they are not mandatory for model's correctness. Constraints (12) are the definition of batch latenesses and force L_{max} to be equal to the maximal lateness. Finally, a solution to the mathematical formulation is a feasible assignment of the jobs to an EDD-sequence of batches. As opposed to the constraint programming formulation, the solution encoding does not consider permutations of batches as equivalent solutions.

6.4.2. Comparison

Our approach applies the value selection heuristic *batch fit* and all filtering rules. Comparison have been done with the resolution of the mathematical formulation stated as an Ilog OPL 6.1.1 model and solved by **llog cplex** 11.2.1. **llog cplex** provides a performance tuning tool which tries to find the best combination of its parameters. The tuning has been done on instances with 10 jobs since it is required to solve them optimally several times. First, the tuning changed the kind of *cuts* to reduce the number of branches explored. A *cut* is a constraint added to a model to restrict non-integer solutions that would otherwise be solutions of the continuous relaxation. Second, the tuning favoured feasibility by assigning a job j to a batch k in the first branch $(x_{jk} = 1)$ instead of the default behavior which forbids the packing $(x_{jk} = 0)$.

Figure 5(a) shows the percentage of instances solved optimally within one hour as a function of the time in seconds. Our constraint programming approach outperforms the mathematical formulation: it solves more instances with solution times that are orders of magnitude lower.

Figure 5(b) compares the quality gap on instances with strictly more than 50 jobs. The time limit of **llog cplex** has been increased until 43200 seconds (12h). Each point represents one instance and its x coordinate is the quality gap that we obtained, whereas



Figure 5: Comparison to the mathematical formulation.

its y coordinate is the quality gap obtained with **llog cplex**. All points located above the line (x = y) are upper bounds improved by our approach. Despite a smaller time limit, the constraint programming approach provides better solution than the mathematical formulation on the vast majority of instances. Furthermore, the difference between the two upper bounds is very tight when the mathematical formulation found the best, whereas it can be significant when the constraint approach did.

6.5. Comparison to a branch-and-price

Branch-and-price integrates branch-and-bound and column generation methods for solving large scale integer programs (Barnhart et al., 1998). Daste et al. (2008a) proposed a branch-and-price algorithm with a master problem where each column is a batch. A solution of the master problem is a feasible sequence of batches. At each iteration, the objective of the subproblem is to find one column (batch) which improves the current solution of the master problem. This pricing subproblem is solved by a greedy algorithm, then, if needed, by an exact enumeration method.

Their algorithm has been tested on instances presented in Section 6.1 with less than 50 jobs and a time limit of 3600 seconds (1h). All their experiments have been conducted on a Pentium 2.66GHz with 1GB of RAM. Detailed results on benchmark instances are not mentioned. 55% and 8% of instances of size n = 20 and n = 50 are solved by the branchand-price procedure within one hour. Although their branch-and-price is more efficient than the mathematical formulation, it is less efficient than our constraint programming which solves optimally all instances of size n = 20 with an average solution time lower than 1 second, and 95% of instances of size n = 50 with an average solution time lower than 100 seconds.

7. Conclusion

We have presented a constraint programming approach to minimize the maximal lateness for a batch processing machine on which a finite number of jobs with non-identical sizes must be scheduled. This approach exploits a new optimization constraint based on a relaxed problem which applies cost-based domain filtering rules and the search is enhanced with dedicated branching heuristics. Computational results demonstrate the positive effect of each component and give better solutions with computation times that are orders of magnitude lower than a mathematical formulation or a branch-and-price based on bin packing and sequencing models.

In further research, we will apply our approach to problems with completion time related measures. In addition, subsequent research topics include the study of parallel batching machines and additional constraints, for instance job release dates that remain incompatible with our approach.

AppendixA. Short description of the pack constraint

In this appendix, we briefly describe our implementation of the global constraint **pack** originally proposed by Shaw (2004) for the bin-packing problem. Let recall that this constraint replaces the channeling constraints between assignment variables ($\forall j \in J, B_j = k \Leftrightarrow j \in J_k$) and enforces the consistency between assignments and loads ($\forall k \in K, \sum_{j \in E_k} s_j = S_k$). Furthermore, **pack** propagates the redundant constraint specifying that the sum of the bin loads is equal to the sum of the item sizes ($\sum_J s_j = \sum_K S_k$). First the constraint performs the propagation of typical model which corresponds to Constraints (6) and (7) of the mathematical formulation (see Section 6.4). In addition, it uses propagation rules incorporating knapsack-based reasoning, as well as a dynamic lower bound on the number of bins required.

The additional constraint propagation rules are based upon treating the simpler problem of packing items into a single bin. That is, given a bin, can we find a subset of items that when packed in the bin would bring the load in a given interval? This problem is a type of knapsack or subset sum problem (Kellerer et al., 2004). Proving that there is no solution to this problem for any bin would mean that search could be pruned. If an item appears in every set of items that can be placed in the bin, we can commit it to the bin. Conversely, if the item never appears in such a set, we can eliminate it as a candidate item. So, if no legal packing can attain a given load, it cannot be a legal load for the bin. Shaw (2004) proposed efficient algorithms which do not depend on item sizes or bin capacity, but which do not in general achieve generalized arc consistency on the subset sum problem as opposed to Trick (2003). Next, they adapt bounding procedures to partial assignments. The idea is to transform a partial assignment into a new bin packing instance and to apply a bounding procedure to this new instance. Our implementation uses the lower bound L_{CCM}^{1D} (Carlier et al., 2007), which dominates the lower bound L_{MV}^{1D} originally used.

Our implementation differs also from the one of Shaw (2004) by the addition of variables S_k and M. Indeed, variables S_k helps to express additional constraints and were maintained internally by Shaw (2004). The variable M represents the objective of the bin packing problem and its value can be restricted by other constraints.

Shaw (2004) demonstrated that this global constraint can cut search by orders of magnitude. Additional comparisons showed that the global constraint coupled with a standard packing algorithm based on *complete decreasing* can significantly outperform Martello & Toth (1990)'s procedure.

AppendixB. Algorithm for cost-based domain filtering of assignments

We recall that a simple algorithm can compute marginal costs from scratch for each possible assignment with an overall complexity of $O(n^3)$. In this section, we describe an O(nm) version of this algorithm (cf. Algorithm 1) based on the incremental computation of marginal costs. The principle consists in exploiting the formula below for quick computation of marginal costs by distinguishing several cases. For instance, the assignment of a job to any empty batch leads to the same marginal cost. In fact, the completion time of a bucket q after the assignment of a job j to a batch $k \in \mathcal{D}(B_j)$ is given by:

$$C_q \left(\mathcal{A} \cap \{B_j \leftarrow k\}\right) - C_q(\mathcal{A}) = \begin{cases} 0 & \delta_q < \min(\max(D_k), d_j) (B.1) \\ \max(p_j, \min(P_k)) & d_j \le \delta_q < \max(D_k) & (B.2) \\ \max(p_j - \min(P_k), 0) & \delta_q \ge \max(D_k) & (B.3) \end{cases}$$

Figure B.6 illustrates the idea behind this formula: the assignment of a job j to a batch k leads to transfer the batch k from a bucket q to one of its predecessors or itself, *i.e.* the bucket $q' \leq q$ such that $\delta_{q'} = \min(d_j, \max(D_k))$. First, the sequence of buckets $1, \ldots, q' - 1$ stays unchanged and so does the completion times (B.1) and its maximal lateness which is dominated by $L(\mathcal{A})$. Indeed, the rule (LF) is applied after the initialization of buckets and before the rule (AF). Then, completion times of buckets $q', \ldots, q-1$ are postponed for the updated batch duration $\max(\min(P_k), p_j)$ and therefore, they are given by (B.2). Last, the duration increase $\max(p_j - \min(P_k), 0)$ of batch k is added to the completion time of buckets q, \ldots, n (B.3). If the batch duration increase is zero, then the maximal lateness of buckets q, \ldots, n is also dominated by $L(\mathcal{A})$

The backward algorithm 1 prunes variables B_j according to the rule (AF). The pruning is illustrated in Figure B.7 by considering the assignment of job 4 for the partial assignment \mathcal{A}_2 (see Figure 1 page 8). After initialization, Loop L1 iterates over buckets in descending order to detect infeasible assignments. The maximal lateness $\Lambda_{\geq q}$ and the maximal duration increase Δ of the sub-sequence of buckets q, \ldots, n are updated using recursive formulas derived from Section 4.1.

New batches of $\mathcal{B}_{=q}$ scheduled at the latest at bucket q that are neither empty nor full are now considered for pruning by the algorithm. Loop L2 updates the maximal duration of these batches, but the removal of infeasible assignments is delegated to Constraints (1). This loop detects simultaneously all infeasible assignments of jobs such that $\max(D_k) \leq d_j$, since their marginal costs depend only on the duration increase of the batch when (B.2) is not defined as illustrated in Figure 7(c).

Loop L3 updates incrementally the maximal lateness $\Lambda_{\geq q}^k$ of the sequence q, \ldots, n without the batch k. Then, Loop L4 inspects only assignments of jobs that would effectively

Figure B.6: Anticipating the impact of an assignment on the buckets.

Algorithm 1: Cost-based domain filtering of assignments.

 $\mathcal{B} = \emptyset$; // Set of batches to prune $\Lambda_{>q}=-\infty$; // Lateness of sub-sequence q,\ldots,n for each $k \in K$ do $\Lambda^k_{\geq q} = -\infty$; // Lateness without batch k of buckets q, \ldots, n /* Try assignments in decreasing order of buckets */ L1 for $q \leftarrow n$ to 1 do $\Lambda_{\geq q} = \max(L_q(\mathcal{A}), \Lambda_{\geq q}); //$ Maximal lateness of buckets q, \ldots, n $\Delta = \max(L_{max}) - \Lambda_{\geq q};$ // Maximal duration increase allowed at bucket q/* Update durations of the new batches to prune at bucket q */ $\mathcal{B}_{=q} = \{k \in K_{=q} \mid \min(P_k) > 0 \land (|\mathcal{D}(P_k)| > 1 \lor |\mathcal{D}(D_k)| > 1)\};\$ foreach $k \in \mathcal{B}_{=q}$ do $\max(P_k) \leftarrow \min(P_k) + \Delta;$ L_2 /* assignments of an available job j to other batches */ for each $k \in \mathcal{B}$ do $\Lambda_{\geq q}^k = \max(L_q(\mathcal{A}), \Lambda_{\geq q}^k);$ L3 forall the $j \in J_{=q}$ such that $|\mathcal{D}(B_j)| > 1$ do L4/* Assignment to non-empty batches */ for each $k \in \mathcal{B}$ do L5if $\delta_q < \min(D_k)$ then $\mathcal{B} = \mathcal{B} \setminus \{k\}$; else if $\Lambda_{\geq q}^k + \max(\min(P_k), p_j) > \max(L_{max})$ then $B_j \not\leftarrow k$; /* Assignment to empty batches */ if $p_j > \Delta$ then $\max(B_j) \leftarrow |K^\star|;$ /* Update data structure */ for each $k \in \mathcal{B}_{=q}$ do $\Lambda_{\geq q}^k = \Lambda_{\geq q} - \min(P_k);$ L6 $\mathcal{B} = \mathcal{B} \cup \mathcal{B}_{=q} ;$

lead to the transfer of the target batch to the bucket q. The internal Loop L5 considers batches to prune discovered in previously visited buckets. If the transfer is infeasible because of due date restrictions, then the batch is eliminated from the set \mathcal{B} of batches to prune. Otherwise, the assignment of job j to the batch k is eliminated if the maximal lateness after its transfer into the bucket q exceeds the current upper bound. In fact, the incremental computation of $\Lambda_{\geq q}^k$ reduces the complexity of (B.2) from linear time to constant time. In Figure 7(d), the transfer of batch 2 from the bucket 2 to 1 does not modify the sequence of batches, whereas the transfer of batch 3 from the bucket 3 to 1 leads to swap batches 2 and 3 in Figure 7(e).

Then, the algorithm simultaneously considers all assignments of a job to an empty batches for which (B.2) and (B.3) are equals, since the duration increase is equal to the updated duration $(\min(P_k) = 0)$ as illustrated in Figure 7(f). In this case, B_j is restricted to values lower than the index $|K^*|$ of the last non-empty batch.

Lastly, Loop L6 initializes the value $\Lambda_{\geq q}^k$ of new batches to prune by subtracting their contribution to the lateness of the sub-sequence. Then, the set \mathcal{B} of batches to prune is updated.

The complexity of Algorithm 1 depends only on its nested loops because instructions are done in constant time. A simple calculation shows that the complexity depends only on nested Loops L1, L4, and L5. In the worst case, nested Loops L1 and L4 iterates over a partition of the jobs J which is done in O(n). Loop L5 is done in O(m), because $\mathcal{B} \subseteq K$ is a subset of batches. Therefore, the complexity of Algorithm 1 is O(nm).

Figure B.7: Impact on the solution of the relaxed problem of the assignment of job 4 to a batch for the partial assignment $\mathcal{A}_2 = \{B_1 \leftarrow 1, B_2 \leftarrow 2, B_3 \leftarrow 3, B_4 \in [1, 4]\}.$

AppendixC. Results summary

Table C.2 page 23 summarizes the best known lower and upper bounds. The first column (#) reports the index of the instance. Columns 2–3 report optimum values (L_{max}) for problem sizes n = 10, 20. Columns 4–9 report lower (*lb*) and upper (*ub*) bounds for problem sizes n = 50, 75, 100 (the lower bound is not given when the upper bound has been proven optimal).

	n = 10	n = 20	n = 50		<i>n</i> =	n = 75		n = 100	
#	Lmax	L_{max}	lb	ub	lb	ub	lb	ub	
1	71	389		1349	1548	1712	2647	2939	
2	59	282		893	1316	1600	1964	2430	
3	37	63		132	-27	263	58	427	
4	-103	-147		20	372	453	72	379	
5	132	337		1502		1481	1866	2218	
6	21	254		1139	1409	1725	1851	2246	
7	-90	33		74	59	149	386	400	
8	-17	-61		11	9	58	66	190	
9	163	309		1043	1807	2164	1863	2575	
10	7	313		1114	1327	1435	2509	2863	
11	50	-39		285		201	98	310	
12	-15	-119		91	314	396	-62	316	
13	58	278		948	1743	2033	2421	2702	
14	4	209		1322	1450	1830	1770	2024	
15	-21	127		178	-31	302	158	633	
10	-49	-20		252	-30	1010	280	0/0	
10	150	303		1022	10598	1492	949E	2840	
10	203	200		267	1052	1423	2400	2090	
20	-101	-16		202	-41	288	-120	420	
20	-45	-10		202	100	200	-120	410	
21	59	240		1156	1754	2054	2440	2592	
22	44	291		1128	150	2105	1468	2167	
23	-66	17		181	159	181	62	130	
24	-43	203		-102	-64	139	473	609	
25	130	020		1974	1405	1502	1065	2609	
20	58	214		169	1047	1092	1905	2400	
21	-105	-34		102	102	273	125	482	
20	140	257		1255	1531	1687	2420	2607	
30	49	353		1010	1528	1714	1829	2438	
31	-13	106	I	269	67	200		338	
32	-64	-33		324	-72	238	10	449	
33	197	423		1117	1776	2032	2134	2528	
34	44	293		936	1132	1421	2156	2605	
35	27	53		292	204	301	-18	496	
36	-32	-55	-41	30	-19	47	-32	259	
37	81	227		1328	1535	1880	2516	2780	
38	10	216		1025	1260	1452	2027	2323	
39	11	58		114	354	355	96	414	
40	-24	14	15	102	204	261	175	393	

Table C.2: Summary of best known lower and upper bounds for the benchmark instances of Daste et al..

Azizoglu, M., & Webster, S. (2000). Scheduling a batch processing machine with non-identical job sizes. International Journal of Production Research, 38, 2173–2184.

Azizoglu, M., & Webster, S. (2001). Scheduling a batch processing machine with incompatible job families. Computers and Industrial Engineering, 39, 325–335.

Baptiste, P., Le Pape, C., & Nuijten, W. (2001). Constraint-Based Scheduling: Applying Constraint Programming to Scheduling Problems. Kluwer.

Barnhart, C., Johnson, E. L., Nemhauser, G. L., Savelsbergh, M. W. P., & Vance, P. H. (1998). Branchand-price: column generation for solving huge integer programs. Operations Research, 46, 316–329.

- Boucher, E., Bachelu, A., Varnier, C., Baptiste, P., & Legeard, B. (1997). Multi-criteria comparison between algorithmic, constraint logic and specific constraint programming on a real scheduling problem. In Proceedings of the 3rd International Conference on the Practical Application of Constraint Technology (pp. 47–64). Blackpool: The Practical Application Ltd.
- Brucker, P., Gladky, A., Hoogeveen, H., Koyalyov, M., Potts, C., Tautenham, T., & van de Velde, S. (1998). Scheduling a batching machine. *Journal of Scheduling*, 1, 31–54.
- Carlier, J., Clautiaux, F., & Moukrim, A. (2007). New reduction procedures and lower bounds for the two-dimensional bin packing problem with fixed orientation. *Computers & Operations Research*, 34, 2223–2250.
- Choco Team (2011). Choco: an open source java constraint programming library. http://choco.mines-nantes.fr.
- Damodaran, P., Manjeshwar, P. K., & Srihari, K. (2006). Minimizing makespan on a batch-processing machine with non-identical job sizes using genetic algorithms. *International Journal of Production Economics*, 103, 882–891.
- Damodaran, P., Srihari, K., & Lam, S. S. (2007). Scheduling a capacitated batch-processing machine to minimize makespan. *Robotics and Computer-Integrated Manufacturing*, 23, 208–216.
- Daste, D., Gueret, C., & Lahlou, C. (2008a). A Branch-and-Price algorithm to minimize the maximum lateness on a batch processing machine. In *Proceedings of the 11th international workshop on Project Management and Scheduling PMS'08* (pp. 64–69). IStanbul Turquie.
- Daste, D., Gueret, C., & Lahlou, C. (2008b). Génération de colonnes pour l'ordonnancement d'une machine à traitement par fournées. In 7ème conférence internationale de modélisation et simulation (MOSIM 2008) (pp. 1783–1790). Paris France volume 3.
- Dupont, L., & Dhaenens-Flipo, C. (2002). Minimizing the makespan on a batch machine with nonidentical job sizes: An exact procedure. Computers & Operations Research, 29, 807–819.
- Felizari, L. C., de Arruda, L. V., Lüders, R., & Stebel, S. L. (2009). Sequencing batches in a real-world pipeline network using constraint programming. In C. A. O. d. N. Rita Maria de Brito Alves, & E. C. Biscaia (Eds.), 10th International Symposium on Process Systems Engineering: Part A (pp. 303 – 308). Elsevier volume 27 of Computer Aided Chemical Engineering.
- Focacci, F., Lodi, A., & Milano, M. (1999). Cost-based domain filtering. In J. Jaffar (Ed.), Proceedings of the 5th International Conference on Principles and Practice of Constraint Programming (CP 1999) (pp. 189–203). Springer volume 1713 of Lecture Notes in Computer Science.
- Fukunaga, A. S., & Korf, R. E. (2007). Bin completion algorithms for multicontainer packing, knapsack, and covering problems. *Journal of Artificial Intelligence Research*, 28, 393–429.
- Garey, M. R., & Johnson, D. S. (1979). Computers and Intractability: A Guide to the Theory of NP-Completeness. New York, USA: W. H. Freeman and Company.
- Gent, I. P., & Walsh, T. (1997). From approximate to optimal solutions: Constructing pruning and propagation rules. In Proceedings of the 15th International Joint Conference on Artificial Intelligence (IJCAI 1997) (pp. 1396–1401).
- Ghazvini, F. J., & Dupont, L. (1998). Minimizing mean flow times criteria on a single batch processing machine with non-identical jobs sizes. *International Journal of Production Economics*, 55, 273–280.

IBM (2011). IBM Ilog Cplex Optimizer. http://www-01.ibm.com/software/integration/ optimization/cplex-optimizer/.

- Ikura, Y., & Gimple, M. (1986). Scheduling algorithms for a single batch processing machine. Operations Research Letters, 5, 61–65.
- Kashan, A. H., Karimi, B., & Ghomi, S. F. (2009). A note on minimizing makespan on a single batch processing machine with non-identical job sizes. *Theoretical Computer Science*, 410, 2754–2758.
- Kellerer, H., Pferschy, U., & Pisinger, D. (2004). Knapsack Problems. Springer, Berlin, Germany.
- Kotecha, P., Gudi, R., Bhushan, M., & Kapadi, M. (2007). On the determination of multiple solutions for batch scheduling using constraint programming. In AIChE Annual Meeting, Conference Proceedings. Lawler, E. (1973). Optimal sequencing of a single machine subject to precedence constraints. Manage-
- *ment Science*, 19, 544–546. Liu, L., Ng, C., & Cheng, T. (2007). Scheduling jobs with agreeable processing times and due dates on
- Liu, L., Ng, C., & Cheng, I. (2007). Scheduling jobs with agreeable processing times and due dates on a single batch processing machine. *Theoretical Computer Science*, 374, 159–169.
- Malve, S., & Uzsoy, R. (2007). A genetic algorithm for minimizing maximum lateness on parallel identical batch processing machines with dynamic job arrivals and incompatible job families. *Computers & Operations Research*, 34, 3016–3028.
- Martello, S., & Toth, P. (1990). Knapsack Problems: Algorithms and Computer Implementations. New York: Wiley.
- Mehta, S. V., & Uzsoy, R. (1998). Minimizing total tardiness on a batch processing machine with

incompatible job families. IIE Transactions, 30, 165-178.

- Parsa, N. R., Karimi, B., & Kashan, A. H. (2010). A branch and price algorithm to minimize makespan on a single batch processing machine with non-identical job sizes. *Computers & Operations Research*, 37, 1720–1730.
- Perez, I., Fowler, J., & Carlyle, W. (2000). Minimizing total weighted tardiness on a single batch process machine with incompatible job families. *Computers & Operations Research*, 32, 327–341.
- Potts, C. N., & Kovalyov, M. Y. (2000). Scheduling with batching: a review. European Journal of Operational Research, 120, 228-249.
- Rossi, F., Beek, P. v., & Walsh, T. (2006). Handbook of Constraint Programming (Foundations of Artificial Intelligence). New York, USA: Elsevier Science Inc.
- Sabouni, M. Y., & Jolai, F. (2010). Optimal methods for batch processing problem with makespan and maximum lateness objectives. Applied Mathematical Modelling, 34, 314–324.
- Shaw, P. (2004). A constraint for bin packing. In M. Wallace (Ed.), Proceedings of the 10th International Conference on Principles and Practice of Constraint Programming (CP 2004) (pp. 648–662). Springer volume 3258 of Lecture Notes in Computer Science.
- Trick, M. A. (2003). A dynamic programming approach for consistency and propagation for knapsack constraints. Annals of Operations Reseach, 118, 73–84.
- Uzsoy, R. (1995). Scheduling batch processing machines with incompatible job families. International Journal of Production Research, 33, 2685–2708.
- Vilím, P. (2007). Global Constraints in Scheduling. Ph.D. thesis Charles University in Prague, Faculty of Mathematics and Physics, Department of Theoretical Computer Science and Mathematical Logic Prague, Czech Republic.
- Wang, C., & Uzsoy, R. (2002). A genetic algorithm to minimize maximum lateness on a batch processing machine. Computers & Operations Research, 29, 1621–1640.
- Webster, S., & Baker, K. (1995). Scheduling groups of jobs on a single machine. Operations Research, 43, 692–703.
- Zeballos, L. J., & Henning, G. P. (2003). A constraint programming approach to the multi-stage batch scheduling problem. In Proceedings of the 4th International Conference on Foundations of Computer-Aided Operations. Coral Springs, Florida, U.S.A.