

Models for Radiation Therapy Patient Scheduling

Sara Frimodig^{1,2} and Christian Schulte¹

¹ KTH Royal Institute of Technology, Sweden, sarhal@kth.se, cschulte@kth.se

² RaySearch Laboratories, Stockholm, Sweden

Abstract. In Europe, around half of all patients diagnosed with cancer are treated with radiation therapy. To reduce waiting times, optimizing the use of linear accelerators for treatment is crucial. This paper introduces an Integer Programming (IP) and two Constraint Programming (CP) models for the non-block radiotherapy patient scheduling problem. Patients are scheduled considering priority, pattern, duration, and start day of their treatment. The models include expected future patient arrivals. Treatment time of the day is included in the models as time windows which enable more realistic objectives and constraints. The models are thoroughly evaluated for multiple different scenarios, altering: planning day, machine availability, arrival rates, patient backlog, and the number of time windows in a day. The results demonstrate that the CP models find feasible solutions earlier, while the IP model reaches optimality considerably faster. All models are shown to be sensitive to an increase in the number of time windows per day.

1 Introduction

Radiation therapy (RT), chemotherapy, and surgery are the most commonly used cancer therapies worldwide. In RT, machines called linear accelerators (LINACs) deliver beams of radiation to the tumor in order to kill malignant tumor cells.

A long waiting time between when a patient is ready for RT and when the treatment starts has a negative effect on its outcome due to for example tumor growth, psychological distress of the patient, and progressed symptoms [8,12,14,22,30]. Hence, many cancer institutes worldwide have adopted waiting time targets that determine the maximum waiting time before treatment starts.

The intent of RT treatments is either curative or palliative, where the latter mainly aims to provide pain relief. Furthermore, cancer patients are generally divided into three urgency levels depending on the site of the cancer, treatment intent, and the size and progress of the tumor. The waiting time targets depend on the patient's urgency level.

RT treatments are generally divided into a number of fractions that are delivered once a day and together sum up to the planned radiation dose. The duration of the fractions vary between patients due to for example treatment technique and tumor complexity [16]. There are also many uncertainties in the RT process, including for example patient inflow and unexpected machine failures.

The scheduling of RT patients on LINACs can be divided into block or non-block systems [10]. Block scheduling systems divide days into slots of equal duration, whereas non-block systems allow for different treatment durations. Block systems are more widely used, but have severe drawbacks since there is no way to control the variability of treatment time, which can generate costs related to machine underutilization, staff overtime, and patient waiting time.

Scheduling patients is mostly done manually and is a considerable challenge for RT clinics. Designing more efficient appointment schedules would be of great significance and could potentially save lives. In order to improve the scheduling of radiotherapy patients, this paper makes the following contributions:

- Two CP models and one IP model of the non-block RT patient scheduling problem are introduced that take expected future patients into account. To the best of our knowledge, these are the first CP models as well as the first IP model to include expected future patient arrivals.
- The treatment time of the day is included in the models which for the first time supports objectives and constraints on treatment time of the day.
- The models capture real-world constraints such as non-consecutive treatment days, different treatment durations, allowed start days, and patient priorities.
- The models are compared using a patient arrival model and several numerical experiments based on data from a European cancer clinic.

Plan of the Paper. Section 2 discusses related work. Section 3 presents the setup of the problem. Section 4 describes the models, followed by a description of the search heuristics in Section 5. The models are experimentally evaluated in Section 6. Section 7 presents conclusions and potential extensions.

2 Related Work

Related work on optimal scheduling in health care has mainly focused on nurse scheduling (see [5]), outpatient assignment (see [7]), and surgery scheduling (see [21]). The RT patient scheduling problem shares some characteristics with these problems, but has particular attributes that make it difficult to apply the models and methods proposed in the literature.

Scheduling of RT patients is a relatively young field with limited literature. In 2016, a review on the literature using operations research for resource planning in RT was published [31]. The authors found 12 papers addressing the problem of scheduling patients on LINACs, where the first ones are published in 2006.

In [18], the authors show that RT patient scheduling can be seen as a special case of a dynamic job-shop problem. They review different exact and metaheuristic methods suitable for solving the problem. A heuristic that schedules patients forward from the first feasible start date (ASAP) is developed in [25]. A local search heuristic that outperforms the ASAP approach is developed in [26].

The first use of IP for optimization of RT appointments is presented in [9] and [10]. Another IP model for non-block scheduling is presented in [17]. A limitation in these papers is that they do not consider all the constraints present

in RT scheduling, such as for example treatments on non-consecutive days and LINAC eligibility. In [6], the authors develop an IP model that includes more realistic constraints, but still using a myopic scheduling policy, i.e., not taking future patient arrivals into account.

Using a block scheduling strategy, [27] presents a method for advance RT patient scheduling, where appointments are scheduled in advance of the service date with future demand still unknown. A Markov decision process (MDP) and approximate dynamic programming are used to solve the problem, and they achieve very good results. In [13], the authors use the same problem setup but also include patient cancellations using simulation-based solution methods. In these papers, there is no way to include time of the day for the treatments.

A hybrid combining stochastic and online optimization is presented in [19]. The authors use a block-scheduling strategy to schedule curative patients at the same time every day and require that patients leave the center with their appointment, which calls for short computation times. This is different from earlier published methods that all schedule multiple patients in a batch.

CP has been used in RT treatment planning [2] and in chemotherapy patient scheduling [15]. Scheduling is a field where CP has shown to be effective, see for example [3]. A comprehensive review of operations research methods for optimization in radiation oncology is presented in [11], where it is stated that CP has not yet had a significant impact on medical physics.

3 Radiation Therapy Patient Scheduling

This section introduces the RT patient scheduling problem, the assumptions made in this paper, and some fundamental modeling aspects.

Time. Radiation therapy clinics have different routines for scheduling patients on LINACs. Some gather patients into a *batch* and schedule them once or several times a day, while others immediately schedule a single patient. This paper focuses on batch scheduling and assumes that the scheduling is done at the end of each day taking patients from previous days into account. As previously stated, RT clinics can be divided into two categories; those who use *block* and *non-block* scheduling systems. In order to be able to control the variability of treatment time, this paper uses a *non-block* scheduling strategy.

A day is divided into *time windows*. A time window is typically 1.5 – 4 hours while a treatment takes 10–45 minutes. Patients are assigned to windows instead of specific start times as this leads to simpler and more efficient models while maintaining an adequate level of detail from a clinical perspective.

Patients. A physician assigns a *priority* to each patient based on urgency and treatment intent (palliative or curative). It is assumed that there are three priority groups, and therefore three waiting time targets: 2 days for priority A (the highest), 14 days for priority B, and 28 days for priority C patients.

A patient is assigned to a *treatment protocol*, which states the *fractionation scheme* (that is, how many days the patient is to be treated and with which

frequency) and the *duration* of each treatment. Different protocols have different allowed *start days*, which enforces that fewer patients are scheduled on weekends. Some protocols also specify that treatment must start on a *certain time of the day*. In this paper, the protocols used are from a large cancer center in Europe.

The scheduled times are communicated to the patient at most one week before the start date or immediately for priority A patients. Fractions are communicated and cannot be re-planned, as this is the collaboration clinic’s approach. The schedule can change until being communicated: booking decisions are postponed to the next day if patients are scheduled more than a week away.

When creating a patient schedule in practice, the booking administrator needs to make sure that there is room in the schedule for more urgent future patients. In most cases, this is done by leaving some empty time on each machine. In the models, the *expected future patient arrivals* are included to predict the expected utilization of resources. Only the expected future patients who have a waiting time target shorter than the maximum waiting time target of current patients are included. This is as patients with longer waiting time targets will have little or no effect on the current schedule.

An overall *arrival rate* can be extracted from historical data for each clinic, as well as the proportion of arrivals for each priority group. This paper uses the same proportions between the priorities as [19]: 31% are priority A, 19% are priority B, and 50% are priority C. The proportions can easily be adjusted to a particular clinic. In the models, a separate priority group D is created for expected future patients of priority A, since the actual priority A patients should have higher priority than expected future priority A patients. These patients are also treated differently in the search heuristics for the CP models, see Section 5. Each arriving patient is randomly assigned to a treatment protocol.

Machines. The radiation is delivered on LINACs. As a rule, larger centers have multiple machine types used for different sorts of treatment and multiple identical machines to have a redundancy in case of machine failures. In small centers, there may be a few identical machines that only serve some treatment types, while more complex cases are sent to larger centers.

In this paper, it is assumed that there are *multiple machines* but only *one machine type*. The machines are exchangeable in that a patient can be scheduled on any machine each day. This scenario is a realistic way of decomposing the multiple machine problem, since the clinics may consider separate scheduling tasks for each machine type. Instead of M machines with W windows each, this is modeled as having one machine with MW windows. Thus, if for example $M = 3$, $W = 4$, then if a patient is scheduled in window 1 – 4, this corresponds to machine 1, window 5 – 8 is machine 2, and window 9 – 12 is machine 3. An alternative would be to model this as one machine with W windows and multiply each window length by M , but this would be a relaxation of the actual problem.

Using multiple machines represents a real-world setting and also allows for having a higher arrival rate. If there were only one machine available, the arrival rate would be very low and dividing very few patients into three different priority groups would not give good statistics for the expected future patient arrivals.

Objective. In this paper, the main objective is to minimize a weighted sum of the violations of the target dates, where the weights reflect that it is worse to violate the target date for a patient with higher priority. The secondary objective is to schedule patients at approximately the same time each day. Some patients may still work or study during treatment and hence prefer mornings or late afternoons. More importantly, the biological effects of the radiation is calculated on having 24 hours between each fraction, however, in most cases it is allowed to deviate from this and it is thus an objective rather than a hard constraint. The second objective is the reason why time windows are used in the models; scheduling a patient in the same time window every day ensures that the treatment is delivered at approximately the same time every day.

4 Models

Three models are developed to capture the RT patient scheduling problem; a scheduling-based CP model, a packing-based CP model and an IP model. These are designed to capture the same real-world constraints and objectives. Using the set $\mathbb{B} := \{0, 1\}$, the inputs to these models are:

$\mathcal{P} = \{1, \dots, P\}$	set of all patients, $P \in \mathbb{N}$
$\mathcal{D} = \{1, \dots, D\}$	set of days in the planning horizon, $D \in \mathbb{N}$
$\mathcal{W} = \{1, \dots, W\}$	set of time windows in a day, $W \in \mathbb{N}$
$w_L = T_s/W$	the window length, where $T_s \in \mathbb{N}$ is the number of time slots during a day with the chosen discretization
$\mathcal{T}_w = \{t_1, t_2, \dots, t_W\}$	set of times when each window starts counting from the beginning of day 1 where $t_1 = 1, \dots, t_i = (i-1)w_L + 1$
$\mathcal{M} = \{1, \dots, M\}$	set of machines, $M \in \mathbb{N}$
$dur_p \in \mathbb{N}$	duration of a fraction for patient p in time slots
$L_p \in \{13, \dots, 47\}$	schedule length for patient p in days
$\mathcal{F} = \{1, \dots, \max(L_p)\}$	set of all treatment days
$FS_p \in \mathbb{B}^{L_p}$	a vector holding the fractionation schedule for patient p , i.e., a vector of ones and zeros representing treatment days and pause days, respectively
$S \in \mathbb{B}^D \times \mathbb{B}^{T_s}$	a matrix holding the partially occupied schedule, where $S_{d,t_s} = 1$ iff time slot t_s on day d is occupied
$\mathcal{A}_p \in \{1, \dots, 7\}$	the set of allowed start days for patient p
$c_p \in \mathbb{N}$	penalty for missing the waiting time target for each priority
$d_{L,p} \in \mathbb{N}$	day limit, i.e., the waiting time target for patient p

4.1 Scheduling-based CP Model

Variables. The basic decision variables of the model are as follows:

$start_{p,f} \in \mathcal{T}_w$	the start time for the window patient $p \in \mathcal{P}$ is scheduled in during treatment day $f \in \mathcal{F}$
$window_{p,f} \in \{0, \dots, W\}$	the window patient $p \in \mathcal{P}$ is scheduled in during treatment day $f \in \mathcal{F}$, where window 0 represents no treatment
$fraction_{p,d} \in \{0, \dots, L_p\}$	fraction that is delivered to patient $p \in \mathcal{P}$ on day $d \in \mathcal{D}$, where fraction 0 represents no treatment

There are also variables that are derived from the basic variables:

$$\begin{array}{ll} day_{p,f} \in \mathcal{D} & \text{day patient } p \in \mathcal{P} \text{ is treated with fraction } f \in \mathcal{F} \\ start_day_p = day_{p,1} & \text{start day for patient } p \in \mathcal{P} \end{array}$$

Constraints. In the scheduling-based CP model, the **cumulative** global constraint [1] is used to ensure that no two treatments overlap. Already scheduled patients, who have a fixed start time and fixed duration in a fixed window, are given as input to the problem and are included in the constraint to ensure that no new patients are scheduled in a window that is already full. The constraint is used “backwards”, setting the duration of each treatment equal to the window length w_L and the resource requirement as the duration of the treatment dur_p .

The variables $day_{p,f}$ and $fraction_{p,d}$ are dual to each other. Their connection is made by the constraints

$$fraction_{p,d} = f, \text{ where } d = day_{p,f} \quad \forall p \in \mathcal{P}, f \in \mathcal{F}, d \in \mathcal{D} \quad (1a)$$

$$(d < start_day_p) \vee (d \geq start_day_p + L_p) \rightarrow fraction_{p,d} = 0 \quad \forall p \in \mathcal{P}, d \in \mathcal{D} \quad (1b)$$

where (1a) is used to express the duality and (1b) is used to express that $fraction_{p,d} = 0$ if patient p is before the start or after the end of treatment.

The day for fraction $f \in \mathcal{F}$ for patient $p \in \mathcal{P}$, $day_{p,f}$, is given by the time when the fraction starts, $start_{p,f}$, divided by the number of time slots T_s , rounded down. In practice, this is expressed using two inequalities in order to avoid performing division on decision variables: $(day_{p,f} - 1)T_s + 1 \leq start_{p,f} \leq day_{p,f}T_s$.

The days are connected to each other by the constraint

$$day_{p,f+1} = day_{p,f} + 1 \quad \forall p \in \mathcal{P}, f \in \mathcal{F}, \quad (2)$$

which means that if treatment day f is on day d , then $f + 1$ is on $d + 1$.

Next, connect $window_{p,f}$ to $start_{p,f}$. The vector FS_p is the fractionation schedule for patient p and is input to the problem given by a protocol. $FS_{p,f} = 1$ corresponds to treatment day f being active for patient p , i.e., treatment is delivered that day. Thus, if the input $FS_{p,f}$ is indeed one, this gives

$$start_{p,f} = (day_{p,f} - 1)T_s + (window_{p,f} - 1)w_L + 1, \quad (3)$$

and for all other $p \in \mathcal{P}, f \in \mathcal{F}$, set $window_{p,f} = 0$. (3) states that the start time for the window patient p is scheduled during fraction f is equal to the start time of that day, plus the start time of the window on that day.

Bounds are also given for when each fraction can start earliest and latest. For example, the patient’s second fraction cannot be on the first day, and similarly, the patient’s first fraction cannot be on the last day:

$$f \leq day_{p,f} \leq D - (L_p - f) \quad \forall f \in \mathcal{F}. \quad (4)$$

Similar constraints limit the start day for each patient to be at latest L_p days from the end of the planning horizon. Patients have to start treatment on an allowed start day, $start_day_p \in \mathcal{A}_p \quad \forall p \in \mathcal{P}$.

To break some dominance, patients of the same priority and treatment protocol are sorted by their waiting time target. A constraint enforces that an earlier target patient always starts their treatment before a later target patient.

Objective Function. The first objective is to start each treatment within the waiting time targets. The target violation is measured as the number of days that the patient misses their treatment target date with

$$target_violation_p = \max(0, start_day_p - d_{L,p}) \quad \forall p \in \mathcal{P}, \quad (5)$$

where $d_{L,p}$ is the day limit for when treatment should start.

The second objective is to schedule patients in the same window each day. Therefore, a penalty is added each time the window is switched. However, since the problem of for example 3 machines with 4 windows is modeled as having 1 machine and 12 windows, there should be no penalty for moving from window 1 to 5 or 9, since they represent the same window but on different machines. Therefore, an extra array m_w is added, which gives the corresponding real machine window for each model window. Only the active treatment days \mathcal{F}_a , when $window_{p,f} \neq 0$, are considered. The constraint is then:

$$m_w = [1, 2, 3, 4, 1, 2, 3, 4, 1, 2, 3, 4] \\ window_diff_p = \sum_{f \in \mathcal{F}_a} (m_w[window_{p,f}] \neq m_w[window_{p,f+1}]) \quad \forall p \in \mathcal{P}. \quad (6)$$

The two parts of the objective function given by (5) and (6) are combined into a weighted sum. Each entry of $window_diff_p$ is taken times 2 is to make it equal to the IP formulation, see Section 4.3. The total objective function is then:

$$\sum_{p \in \mathcal{P}} (100c_p target_violation_p + 2window_diff_p), \quad (7)$$

where c_p are weights for each priority group used to capture that it is worse to violate the target for higher prioritized patients. The weight 100 is to reflect that the waiting time targets are more important than minimizing window switches.

4.2 Packing-based CP Model

In the packing-based CP model, $fraction_{p,d} = 0$ if patient p has not started or has finished treatment on day d is expressed by the **regular** constraint [23], replacing (1b). For example, if the treatment length is 10 days, the regular expression is $0^* \cdot 1 \cdot 2 \cdot \dots \cdot 8 \cdot 9 \cdot 10 \cdot 0^*$ (where \cdot is concatenation) in:

$$\mathbf{regular}([fraction_{p,d} | d \in \mathcal{D}], r) \quad \forall p \in \mathcal{P}. \quad (8)$$

The variables $start_{p,d}$ are removed since the **cumulative** constraint is not used, and constraint (3) is also removed. Instead, regular expressions on the time windows are used to define the treatment patterns, where window 0 corresponds

to no treatment being delivered. The example regular expression $r = 0 \cdot ([1 - 4]^5 \cdot 0^2)^6$ states that $f = 0$ gives $window_{p,0} = 0$, then $window_{p,f} \in \{1, \dots, 4\}$ for five days, followed by two pause days where $window_{p,f} = 0$, and then repeating this pattern six times, leading to the constraint:

$$\text{regular}([window_{p,f} | f \in \{0, \dots, L_p\}], r) \quad \forall p \in \mathcal{P} \quad (9)$$

The global `bin_packing` constraint [29] is used to ensure that the patients fit in each window:

$$\text{bin_packing}([\infty, w_L, w_L, w_L, w_L], [window_{p,f} | p \in \mathcal{P}], [dur_p | p \in \mathcal{P}]) \quad \forall d \in \mathcal{D}, \quad (10)$$

where in (10), constraint (1a) is used to connect d to f . Constraint (10) states that the capacity of window 0 is infinite and the other windows have capacity w_L , i.e., the window length. In practice, the partially occupied input schedule is also taken into account, included as patients that have fixed windows.

4.3 IP Model

Variables. The basic variables of the IP model are (again using $\mathbb{B} = \{0, 1\}$):

$$\begin{aligned} s_{p,d} &\in \mathbb{B} & s_{p,d} = 1 \text{ iff patient } p \in \mathcal{P} \text{ starts treatment on day } d \in \mathcal{D} \\ q_{p,d,f} &\in \mathbb{B} & q_{p,d,f} = 1 \text{ iff patient } p \in \mathcal{P} \text{ has their } f\text{-th treatment day on } d \in \mathcal{D}, f \in \mathcal{F} \\ x_{p,d,w} &\in \mathbb{B} & x_{p,d,w} = 1 \text{ iff patient } p \in \mathcal{P} \text{ is scheduled in window } w \in \mathcal{W} \text{ on day } d \in \mathcal{D} \end{aligned}$$

There are also variables that are derived from the basic variables:

$$\begin{aligned} u_{p,d} &\in \mathbb{B} & u_{p,d} = 1 \text{ iff patient } p \in \mathcal{P} \text{ has an active treatment day on } d \in \mathcal{D} \\ y_{p,d,w} &\in \mathbb{B} & \text{help variable for the objective function where } y_{p,d,w} = 1 \text{ if patient } p \in \mathcal{P} \\ & & \text{is scheduled in window } w \in \mathcal{W} \text{ on day } d \in \mathcal{D}, \text{ and if patient } p \text{ is not} \\ & & \text{scheduled in window } w \text{ on day } d, y_{p,d,w} \in \mathbb{B} \end{aligned}$$

Constraints. Patient $p \in \mathcal{P}$ will be treated for L_p days. Thus, the last day to start treatment is $daylimit_p = D - L_p + 1$. The treatment should start exactly one time on an allowed start day given by \mathcal{A}_p :

$$\sum_{d \subseteq \mathcal{A}_p}^{daylimit_p} s_{p,d} = 1, \quad \forall p \in \mathcal{P}. \quad (11)$$

For the CP models, constraint (4) states which fractions that can be delivered on which days. In the IP model, this is enforced by setting $q_{p,d,f} = 0$ for all $p \in \mathcal{P}, d \in \mathcal{D}, f \notin \mathcal{F}_{p,d}$ where $\mathcal{F}_{p,d} := \{\max(0, d - (D - L_p)), \dots, \min(d, L_p)\}$. Using $\hat{\mathcal{F}}_{p,d}$ to denote $\mathcal{F}_{p,d}$ with the last element excluded, the constraint

$$q_{p,d,f} = q_{p,d+1,f+1} \quad \forall p \in \mathcal{P}, d \in \{1, \dots, D-1\}, f \in \hat{\mathcal{F}}_{p,d} \quad (12)$$

is formulated so that all treatment days are scheduled after each other.

Constraints (13) and (14) are defined to state that the f -th treatment day can only happen once and that patient p is scheduled at most once every day d .

In the CP models, this is enforced by constraints (3) or (8).

$$\sum_{d \in \mathcal{D}} q_{p,d,f} \leq 1 \quad \forall p \in \mathcal{P}, f \in \mathcal{F}_{p,d} \quad (13)$$

$$\sum_{f \in \mathcal{F}_{p,d}} q_{p,d,f} \leq 1 \quad \forall p \in \mathcal{P}, d \in \mathcal{D}. \quad (14)$$

The first treatment day $f = 1$ is given on the start day for each patient:

$$q_{p,d,1} = s_{p,d} \quad \forall p \in \mathcal{P}, d \in \mathcal{D}. \quad (15)$$

In the fractionation schedule for patient p , an active treatment day f gives $FS_{p,f} = 1$. A variable $u_{p,d}$ is introduced so that $u_{p,d} = 1$ iff patient p is during treatment on day d ($q_{p,d,f} = 1$) and has an active treatment day ($FS_{p,f} = 1$) and zero otherwise, thus, it controls if d is an active day or not for patient p :

$$u_{p,d} = \sum_{f \in \mathcal{F}_{p,d}} (q_{p,d,f} FS_{p,f}) \quad \forall p \in \mathcal{P}, d \in \mathcal{D}. \quad (16)$$

Each patient is scheduled in exactly one time window on active treatment days, and not in any window on off-days:

$$\sum_{w \in \mathcal{W}} x_{p,d,w} = u_{p,d} \quad \forall p \in \mathcal{P}, d \in \mathcal{D}. \quad (17)$$

In the CP models, constraints (3) or (9) are used to express the same as (17).

In order to make sure that all treatments fit within each time window, $u_{p,d}$ is used to keep track of if the patient has an active day or not, together with $x_{p,d,w}$, which is one iff patient p is scheduled in window w on day d :

$$\sum_{p \in \mathcal{P}} x_{p,d,w} u_{p,d} dur_p + \sum_{t_s \subseteq w} S_{d,t_s} \leq w_L \quad \forall d \in \mathcal{D}, w \in \mathcal{W}. \quad (18)$$

S is the input schedule, where an element is one iff time slot t_s on day d is occupied. Thus, the sum of the duration of all patients in window w plus the previously occupied slots in that window must be less than or equal the window length w_L . In the CP models, this is enforced by the cumulative constraint (scheduling-based) or constraint (10) (packing-based). A major difference is that the number of constraints in the IP model will grow with the number of time windows, which is not the case in the CP models.

Objective Function. The objective is the same as in the CP model. A penalty is added for the time by which the target is missed:

$$f_{1,p} = \sum_{d=d_{L,p}}^D s_{p,d}(d - d_{L,p}) \quad \forall p \in \mathcal{P}, \quad (19)$$

where $d_{L,p}$ corresponds to the waiting target in days for patient p and $s_{p,d} = 1$ on the start day. (19) corresponds to the CP objective function (5).

The other objective is to schedule the patient in the same time window each day. To do this, a help variable $y_{p,d,w} \in \mathbb{B}$ is introduced so that $y_{p,d,w}$ is one when $x_{p,d,w}$ is one, and the sum of $y_{p,d,w}$'s is one on all days:

$$\begin{aligned} y_{p,d,w} &\geq x_{p,d,w} & \forall p \in \mathcal{P}, d \in \mathcal{D}, w \in \mathcal{W} \\ \sum_{w \in \mathcal{W}} y_{p,d,w} &= 1 & \forall p \in \mathcal{P}, d \in \mathcal{D}, \end{aligned} \quad (20)$$

As for the CP models, the problem with 3 machines with 4 windows each is modeled as having 1 machine with 12 windows, and hence there should be no penalty for switching from window 1 to 5 or 9. Introduce $\mathcal{W}_m = \{1, \dots, W_m\}$ where W_m is the number of windows on each machine. Each window switch is penalized in the second objective function, where we here assume 3 machines:

$$f_{2,p} = \sum_{d \in \hat{\mathcal{D}}} \sum_{w \in \mathcal{W}_m} \left| \sum_{i=\{w, w+W_m, w+2W_m\}} (y_{p,d,i} - y_{p,d+1,i}) \right| \quad \forall p \in \mathcal{P}. \quad (21)$$

Since (21) sums both $0 - 1$ and $1 - 0$ in each switch, there is a penalty of 2 in every switch. To make this equivalent to the CP objective function (6), the factor 2 difference is adjusted for in (7).

A new variable is introduced to avoid absolute value in the objective function, since it makes the model nonlinear: $z_{p,d,w} = \sum_{i=\{w, w+W_m, w+2W_m\}} (y_{p,d,i} - y_{p,d+1,i})$ for $w \in \mathcal{W}_m$. $z_{p,d,w}$ is divided into a positive and negative part; $z_{p,d,w} = z_{p,d,w}^+ - z_{p,d,w}^-$. This then gives

$$\begin{aligned} f_{2,p} &= \sum_{d \in \hat{\mathcal{D}}} \sum_{w \in \mathcal{W}_m} z_{p,d,w}^+ + z_{p,d,w}^- & \forall p \in \mathcal{P} \\ z_{p,d,w}^+ &\geq 0, z_{p,d,w}^- \geq 0 & \forall p \in \mathcal{P}, d \in \mathcal{D}, w \in \mathcal{W}_m. \end{aligned} \quad (22)$$

The two formulations (21) and (22) are equivalent in a minimization setting.

In total, the objective function is equivalent to the CP objective function (7):

$$\sum_{p \in \mathcal{P}} (100c_p f_{1,p} + f_{2,p}), \quad (23)$$

where c_p are weights used to capture that it is worse to violate the target for higher prioritized patients.

5 CP Search

When solving the CP models, the search is conducted in the following order:

1. Assign all priority A patients randomly to a start day as early as possible.

2. Assign all priority C and D patients a start day randomly where the $start_day_p$ variable with the smallest domain size over weighted degree [4] is chosen.
3. Assign the priority B patients to their earliest possible start day.
4. Assign the number of window switches $window_diff_p$ as small as possible for all patients sorted by their duration, with longest duration assigned first.

For easy cases, with few patients to schedule, it is possible to construct deterministic search heuristics for the CP models that perform much better than the random search strategies described above. However, on more difficult cases these search heuristics fail to even find a solution within a reasonable time frame. This is the reason for including randomization in the search process.

When running a minimization problem in CP, a branch-and-bound tree search is used, which follows the same branch until the branch has failed. However, using a restart strategy will cause it to restart from the top node whenever it finds a solution or after a specific number of failures defined by the restart strategy. In this case, using a restart strategy yields better results because the problem is somewhat under-constrained; it is easy to find feasible solutions and relatively few failures occur. Therefore, the Luby restart strategy multiplied by a factor of 100 is used [20]. The interval is a result of testing many different intervals on different problem setups. Doing restarts this often can be compared to approximating a Large Neighborhood Search (LNS) [28] (see also Section 7).

6 Experiments and Results

Multiple different experiments are run to compare and evaluate the three models. The experiments are run on a Windows 10 machine with an Intel Core i9-7940X X-series processor and 64.0 GB of RAM. The IP model is solved using the MIP solver of CPLEX 12.8 in the Python API with default parameters. The CP models use MiniZinc 2.2.2 and are solved with Gecode 6.1.0. Other solvers were tested, such as the lazy clause solver Chuffed, but Gecode gave the best overall results on the tested problem instances.

A simulation engine is built with Python 3.6. In this engine, the first day starts from an empty schedule and patients to be scheduled are assumed to arrive according to a Poisson process. For each simulated day, a patient schedule is created using the previously described models, and the patients are fixed to the schedule if they have a start day within a week (since this is the limit assumed to communicate the schedules to the patients). The schedule from the previous day is used as input for the next day in the simulation, together with the backlog of yet unscheduled patients. For these patients, the waiting time target is adjusted by one day since it is counted from the day the patient is ready for treatment.

The simulation engine is used to generate problem benchmarks, that each have a partially occupied schedule, a patient backlog, and the number of expected patients arriving per day (as discussed in Section 3). Occupation in a schedule is measured as the average occupation of the first two weeks in that schedule. This is not a perfect measure; if the first week is completely booked and the second is

Table 1. Setup of the benchmarks, 3 machines

Benchmark	Occupation (%) (4/6 windows)	Expected number of arriving patients	Number of patients in backlog (4/6 windows)	Number of patients including future arrivals (4/6 windows)	Planning day
a-1	18.5 / 10.2	4	5 / 8	49 / 52	Sunday
a-2	18.5 / 10.2	6	5 / 8	77 / 80	Sunday
a-3	18.5 / 10.2	8	5 / 8	93 / 96	Sunday
a-4	53.9 / 55.0	4	8 / 8	52 / 52	Tuesday
a-5	53.9 / 55.0	6	8 / 8	80 / 80	Tuesday
a-6	53.9 / 55.0	8	8 / 8	96 / 96	Tuesday
w-1	53.9 / 55.0	5	8 / 8	79 / 79	Monday
w-2	53.9 / 55.0	5	8 / 8	79 / 79	Tuesday
w-3	53.9 / 55.0	5	8 / 8	79 / 79	Wednesday
w-4	53.9 / 55.0	5	8 / 8	79 / 79	Thursday
w-5	53.9 / 55.0	5	8 / 8	79 / 79	Friday
w-6	53.9 / 55.0	5	8 / 8	79 / 79	Saturday
w-7	53.9 / 55.0	5	8 / 8	79 / 79	Sunday
l-1	53.9 / 55.0	5	8 / 8	79 / 79	Sunday
l-2	60.9 / 61.8	5	16 / 14	87 / 85	Sunday
l-3	69.8 / 64.8	5	24 / 34	90 / 100	Friday
l-4	73.3 / 69.5	5	45 / 48	111 / 114	Thursday

completely free, some urgent patients will not be able to meet their target dates, although occupancy is 50% in total for these weeks.

Scheduling patients on three different machines, 16 different benchmarks that are grouped into three categories are summarized in Table 1. The categories capture the following aspects:

- a) The average number of patients arriving each day.
- w) Which wweekday the schedule is created.
- l) The load in the input, i.e., the amount of partial occupation in the input schedule and the size of the backlog. These are closely related, since an almost empty schedule does not come with a large backlog, and vice versa. Benchmarks l-3 and l-4 examine scalability and do not represent realistic scenarios, as patients would be transferred to other clinics.

The weights in the objective function are chosen as $c_1 = 10$, $c_2 = 3$, $c_3 = 1$, corresponding to priority group A to C, and is the same as in [24]. Priority D is for expected future patients of priority A and has weight $c_4 = 5$. The day is divided into 4 or 6 time windows. The timeout is set to 6 hours.

The performance of the models is measured as the objective function value as a function of runtime for the benchmarks, for both 4 and 6 windows.

Average Patient Arrival. The results for benchmarks a-1 to a-6 is shown in Figure 1. The results show that except for the 6 window case in benchmark a-4, the IP model reaches optimum considerably faster, and the packing-based CP model outperforms the scheduling-based CP model on all instances. When the arrival rate is low, as in benchmarks a-1 and a-4, the IP and CP models have similar performance. The CP models are however more sensitive to an increase in the average number of patients arriving each day, while the time to reach

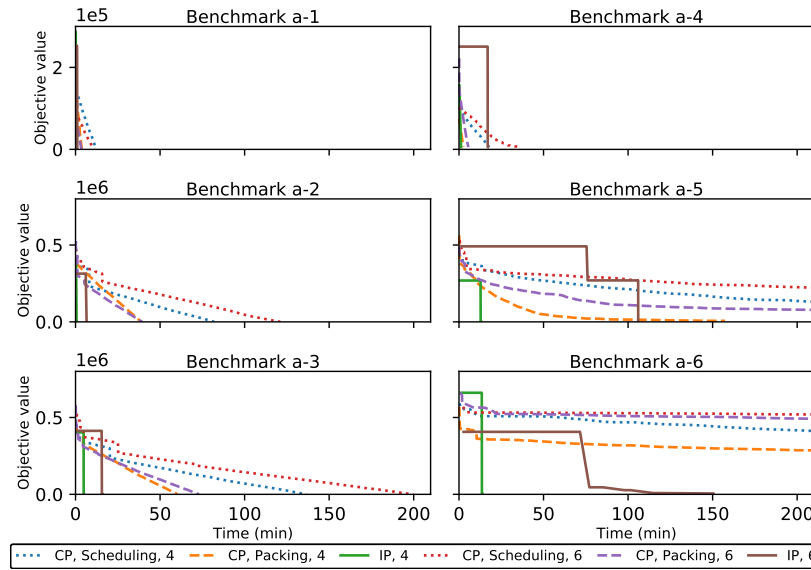


Fig. 1. Results with varying arrival rate, 4 windows and 6 windows

optimality in the 4 window IP model does not change significantly. Increasing the number of time windows from 4 to 6 makes the IP model slower in all cases. When the partial occupation is low (benchmarks a-1 to a-3), the packing-based CP model is not slower with more windows, which is not the case when the partial occupation is higher (benchmarks a-4 to a-6).

Weekday. The results when varying the weekday the schedule is created are presented in Figure 2. Some treatment protocols state that treatment can be initiated only on certain days, but this does not have a large effect on the runtime. An observation is that for the 6-window case, although the IP model reaches optimality faster, the packing-based CP model has a better objective value initially. Benchmarks w-1 to w-7 show that if the time limit is short, the packing-based CP model performs better than the IP model.

Patient Load. Altering the load, the results can be seen in Figure 3. Again, the IP model outperforms the CP models in finding optimality. In benchmarks l-3 and l-4, in 6 hours of runtime the scheduling-based CP model does not find any solutions, the packing-based CP model does not find a solution in the 6-window case, and the IP model does not reach optimality in the 6-window case. However, these cases represent too heavy a load to be realistic scenarios.

Time to Feasible Solution. Figures 1 to 3 show that the IP model outperforms the CP models in finding an optimal solution. However, the packing-based CP model is in almost all cases faster to find a feasible solution, see Figure 4.

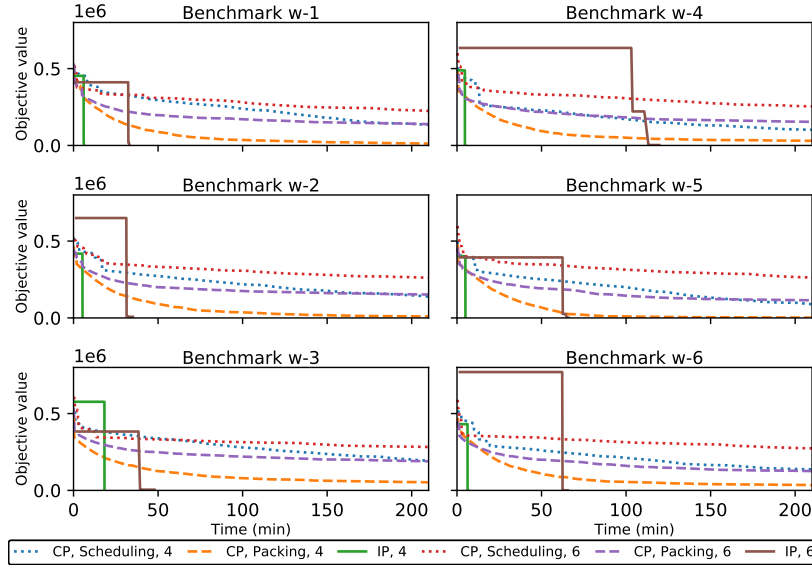


Fig. 2. Results with varying planning day, 4 windows and 6 windows

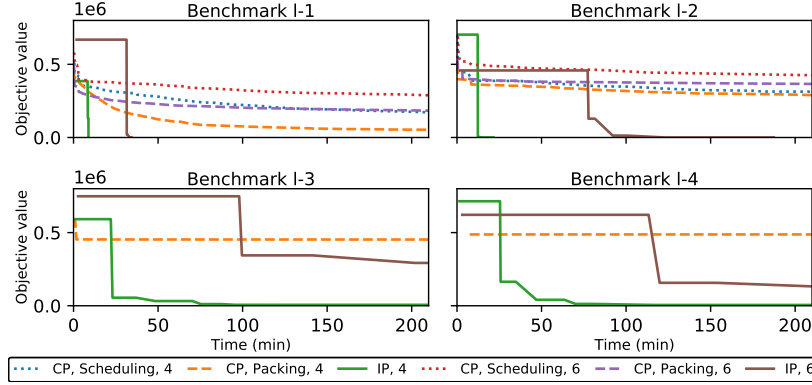


Fig. 3. Results when varying the load, 4 windows and 6 windows

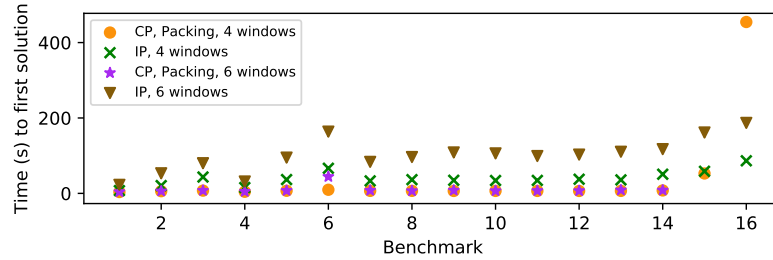


Fig. 4. Time to find the first feasible solution for each of the benchmarks

7 Conclusions and Future Work

Radiation therapy is one of the most commonly used cancer therapies worldwide. The waiting times can be reduced by the optimization of the use of LINACs for treatment. This paper introduces three different models that capture the non-block radiotherapy patient scheduling problem; one IP model, a scheduling-based CP model, and a packing-based CP model. The patients have different priority levels, treatment patterns, treatment duration, and start days for treatment. The expected future patient arrivals are included in the models to predict future resource utilization. The models are evaluated for multiple different scenarios. The results show that the packing-based CP model outperforms the scheduling-based CP model on all problem instances. In general, the CP models find feasible solutions faster than the IP model, however, the IP model reaches optimality considerably faster. All models are sensitive to an increase in the number of time windows per day.

Future Work. To make the models more realistic in the future, overtime on the machines and cancellation of treatments should be included in the models. Cancellation of treatments is common in RT, and allowing for overtime on the machines is important since it is often used in practice to reduce waiting times. Both aspects will increase the complexity of the models, and overtime will most likely make the IP model nonlinear. Another future direction is to extend the models to include multiple machine types, which also increases their complexity.

For the models to be able to handle these complexities, a potential future extension is to take advantage of the strengths of each model in an IP/CP hybrid. Another option would be to use some decomposition method on the IP model, for example Benders decomposition. To better explore the search space, an option is to use Large Neighborhood Search (LNS) when solving the CP models.

A potential improvement of the stochastic aspect of the models is to use scenario-based probabilities instead of expected values when accounting for patients arriving in the future.

The models have so far been compared with each other. The next step would be to see how the models perform over time, using the simulation engine and comparing to both a myopic scheduling strategy (not taking future patient arrivals into account) and to historical data from a clinic.

Acknowledgments. The authors thank Per Enqvist at KTH and Kjell Eriksson at RaySearch Laboratories for a fruitful collaboration. The authors are grateful for insightful discussions about the CP models with Mats Carlsson and Peter J. Stuckey and constructive comments from the anonymous reviewers.

References

1. Aggoun, A., Beldiceanu, N.: Extending Chip in order to solve complex scheduling and placement problems. *Mathematical and Computer Modelling* **17**(7), 57 – 73 (1993). [https://doi.org/10.1016/0895-7177\(93\)90068-A](https://doi.org/10.1016/0895-7177(93)90068-A)

2. Baatar, D., Boland, N., Brand, S., Stuckey, P.J.: CP and IP approaches to cancer radiotherapy delivery optimization. *Constraints* **16**, 173–194 (2011). <https://doi.org/10.1007/s10601-010-9104-1>
3. Barták, R., Salido, M., Rossi, F.: New trends in constraint satisfaction, planning, and scheduling: A survey. *The Knowledge Engineering Review* **25**, 249–279 (09 2010). <https://doi.org/10.1017/S0269888910000202>
4. Boussemart, F., Hemery, F., Lecoutre, C., Sais, L.: Boosting systematic search by weighting constraints. In: de Mántaras, R.L., Saitta, L. (eds.) *Sixteenth European Conference on Artificial Intelligence*. pp. 146–150. IOS Press, Valencia, Spain (Aug 2004)
5. Burke, E.K., De Causmaecker, P., Berghe, G.V., Van Landeghem, H.: The State of the Art of Nurse Rostering. *Journal of Scheduling* **7**(6), 441–499 (nov 2004). <https://doi.org/10.1023/B:JOSH.0000046076.75950.0b>
6. Burke, E.K., Leite-Rocha, P., Petrovic, S.: An integer linear programming model for the radiotherapy treatment scheduling problem. *arXiv e-prints arXiv:1103.3391* (2011)
7. Cayirli, T., Veral, E.: Outpatient Scheduling in Health Care: a Review of Literature. *Production and Operations Management* **12**(4), 519–549 (2009). <https://doi.org/10.1111/j.1937-5956.2003.tb00218.x>
8. Chen, Z., King, W., Pearcey, R., Kerba, M., Mackillop, W.J.: The relationship between waiting time for radiotherapy and clinical outcomes: A systematic review of the literature. *Radiotherapy and Oncology* **87**(1), 3 – 16 (2008). <https://doi.org/10.1016/j.radonc.2007.11.016>
9. Conforti, D., Guerriero, F., Guido, R.: Optimization models for radiotherapy patient scheduling. *4OR* **6**(3), 263–278 (Sep 2008). <https://doi.org/10.1007/s10288-007-0050-8>
10. Conforti, D., Guerriero, F., Guido, R.: Non-block scheduling with priority for radiotherapy treatments. *European Journal of Operational Research* **201**(1), 289–296 (2010). <https://doi.org/10.1016/j.ejor.2009.02.016>
11. Ehrgott, M., Holder, A.: Operations Research Methods for Optimization in Radiation Oncology. *Journal of Radiation Oncology Informatics* **6**(1), 1–41 (2014). <https://doi.org/10.5166/jroi-6-1-21>
12. Fortin, A., Bairati, I., Albert, M., Moore, L., Allard, J., Couture, C.: Effect of treatment delay on outcome of patients with early-stage head-and-neck carcinoma receiving radical radiotherapy. *International Journal of Radiation Oncology Biology Physics* **52**(4), 929–936 (2002). [https://doi.org/10.1016/S0360-3016\(01\)02606-2](https://doi.org/10.1016/S0360-3016(01)02606-2)
13. Gocgun, Y.: Simulation-based approximate policy iteration for dynamic patient scheduling for radiation therapy. *Health Care Management Science* **21**(3), 317–325 (2018). <https://doi.org/10.1007/s10729-016-9388-9>
14. Gomez, D.R., Liao, K.P., Swisher, S.G., Blumenschein, G.R., Erasmus, J.J., Buchholz, T.A., Giordano, S.H., Smith, B.D.: Time to treatment as a quality metric in lung cancer: Staging studies, time to treatment, and patient survival. *Radiotherapy and Oncology* **115**(2), 257–263 (2015). <https://doi.org/10.1016/j.radonc.2015.04.010>
15. Hahn-Goldberg, S., Beck, J.C., Carter, M.W., Trudeau, M., Sousa, P., Beattie, K.: Solving the chemotherapy outpatient scheduling problem with constraint programming. *Journal of Applied Operational Research* **6**(3), 135–144 (2014)
16. Halperin, E.C., Wazer, D.E., Brady, L.W., Perez, C.A.: *Perez and Brady’s principles and practice of radiation oncology*. Lippincott Williams and Wilkins, sixth edn. (2013)

17. Jacquemin, Y., Marcon, E., Pommier, P.: A pattern-based approach of radiotherapy scheduling. In: IFAC Proceedings Volumes. vol. 44, pp. 6945–6950 (jan 2011). <https://doi.org/10.3182/20110828-6-IT-1002.00502>
18. Kapamara, T., Sheibani, K., OCL, H., Reeves, C., Petrovic, D.: A review of scheduling problems in radiotherapy. In: Proceedings of the International Control Systems Engineering Conference on Systems Engineering (ICSE2006). pp. 201–207 (2006)
19. Legrain, A., Fortin, M.A., Lahrichi, N., Rousseau, L.M.: Online stochastic optimization of radiotherapy patient scheduling. *Health Care Management Sciences* (18), 110–123 (2015). <https://doi.org/10.1007/s10729-014-9270-6>
20. Luby, M., Sinclair, A., Zuckerman, D.: Optimal speedup of Las Vegas algorithms. *Information Processing Letters* **47**, 173–180 (1993)
21. May, J.H., Spangler, W.E., Strum, D.P., Vargas, L.G.: The Surgical Scheduling Problem: Current Research and Future Opportunities. *Production and Operations Management* **20**, 392–405 (2011). <https://doi.org/10.1111/j.1937-5956.2011.01221.x>
22. O’Rourke, N., Edwards, R.: Lung cancer treatment waiting times and tumour growth. *Clinical Oncology* **12**(3), 141–144 (2000). <https://doi.org/10.1053/clon.2000.9139>
23. Pesant, G.: A regular language membership constraint for finite sequences of variables. In: Wallace [32], pp. 482–495. <https://doi.org/10.1007/b100482>
24. Petrovic, D., Castro, E., Petrovic, S., Kapamara, T.: Radiotherapy Scheduling. In: Uyar, A., Ozcan, E., Urquhart, N. (eds.) *Automated Scheduling and Planning, Studies in Computational Intelligence*, vol. 505, pp. 155–189. Springer, Berlin, Heidelberg (2013). <https://doi.org/10.1007/978-3-642-39304-4>
25. Petrovic, S., Leung, W., Song, X., Sundar, S.: Algorithms for radiotherapy treatment booking. In: Qu, R. (ed.) *25th Workshop of the UK Planning and Scheduling Special Interest Group (PlanSIG2006)*. pp. 105–112 (04 2006)
26. Riff, M.C., Cares, J.P., Neveu, B.: RASON: A new approach to the scheduling radiotherapy problem that considers the current waiting times. *Expert Systems with Applications* **64**, 287–295 (dec 2016). <https://doi.org/10.1016/J.ESWA.2016.07.045>
27. Sauré, A., Patrick, J., Tyldesley, S., Puterman, M.L.: Dynamic multi-appointment patient scheduling for radiation therapy. *European Journal of Operational Research* **223**(2), 573 – 584 (2012). <https://doi.org/10.1016/j.ejor.2012.06.046>
28. Shaw, P.: Using constraint programming and local search methods to solve vehicle routing problems. In: Maher, M., Puget, J.F. (eds.) *Fourth International Conference on Principles and Practice of Constraint Programming. LNCS*, vol. 1520, pp. 417–431. Springer Berlin Heidelberg, Berlin, Heidelberg (1998). https://doi.org/10.1007/3-540-49481-2_30
29. Shaw, P.: A constraint for Bin Packing. In: Wallace [32], pp. 648–662. <https://doi.org/10.1007/b100482>
30. Van Harten, M.C., Hoebbers, F.J., Kross, K.W., Van Werkhoven, E.D., Van Den Brekel, M.W., Van Dijk, B.A.: Determinants of treatment waiting times for head and neck cancer in the Netherlands and their relation to survival. *Oral Oncology* **51**(3), 272–278 (2015). <https://doi.org/10.1016/j.oraloncology.2014.12.003>
31. Vieira, B., Hans, E.W., Van Vliet-Vroegindewey, C., Van De Kamer, J., Van Harten, W.: Operations research for resource planning and -use in radiotherapy: a literature review. *BMC Medical Informatics and Decision Making* **16**(149) (2016). <https://doi.org/10.1186/s12911-016-0390-4>

32. Wallace, M. (ed.): Tenth International Conference on Principles and Practice of Constraint Programming, LNCS, vol. 3258. Springer Berlin Heidelberg, Toronto, Canada (sep 2004). <https://doi.org/10.1007/b100482>