

Fairness over Time in Dynamic Resource Allocation with an Application in Healthcare

ANDREA LODI¹, PHILIPPE OLIVIER², GILLES PESANT², AND SRIRAM SANKARANARAYANAN³

¹*Jacobs Technion-Cornell Institute, Cornell Tech and Technion - IIT, New York, USA*

²*Polytechnique Montréal, Montréal, Canada*

³*IIM Ahmedabad, Ahmedabad, India*

Abstract

Decision making problems are typically concerned with maximizing efficiency. In contrast, we address problems where there are multiple stakeholders and a centralized decision maker who is obliged to decide in a fair manner. Different decisions give different utility to each stakeholder. In cases where these decisions are made repeatedly, we provide efficient mathematical programming formulations to identify both the maximum fairness possible and the decisions that improve fairness over time, for reasonable metrics of fairness. We apply this framework to the problem of ambulance allocation, where decisions in consecutive rounds are constrained. With this additional complexity, we prove structural results on identifying fair feasible allocation policies and provide a hybrid algorithm with column generation and constraint programming-based solution techniques for this class of problems. Computational experiments show that our method can solve these problems orders of magnitude faster than a naive approach.

1 Introduction

A resource allocation problem can be defined as the problem of choosing an allocation from a set of feasible allocations to a finite set of stakeholders to minimize some objective function.

Generally, this objective function represents the efficiency of an allocation by considering, for example, the total cost incurred by all the stakeholders (Toshihide Ibaraki, 1988).

However, in some cases, such an objective may be inappropriate due to ethical or moral considerations. For instance, to whom should limited medical supplies go in a time of crisis (Heier Stamm et al., 2017)? These situations call for a different paradigm of decision-making that addresses the need of *fairness* in resource allocation.

Resource allocation problems trace their roots to the 1950s when the division of effort between two tasks was studied (Koopman, 1953). Fairness in resource allocation problems has been studied since the 1960s. A mathematical model for the fair apportionment of the U.S. House of Representatives considered the minimization of the difference between the disparity of representation between any pairwise states (Burt and Harris, 1963). This problem was revisited about two decades later (Katoh et al., 1985, Zeitlin, 1981). The so-called *minimax* and *maximin* objectives in resource allocation, which achieve some fairness by either minimizing the maximum (or maximizing the minimum) number of resources allocated to entities, have been studied at least since the 1970s (Jacobsen, 1971, Porteus and Yormark, 1972). A recent book on equitable resource allocation presents various models and advocates a lexicographic minimax/maximin approach for fair and efficient solutions (Luss, 2012). This work further discusses multiperiod equitable resource allocation. Fair resource allocation over time is also considered in the dynamic case where resource availability changes over time and is solved using approximation algorithms (Bampis et al., 2018). A combination of efficiency and fairness in resource allocation is common in communication networks (Ogryczak et al., 2014,

Ogryczak, Włodzimierz, 2014), and can be found in various other fields, such as in some pickup and delivery problems (Eisenhandler and Tzur, 2019).

Further, theoretical computer science has also been interested in a range of fair allocation problems. Procaccia (2015) provides a summary of some of the principles of fair allocation of a divisible good. Fair allocation of indivisible good is an active area of research since the seminal paper of Alkan et al. (1991). This literature consider multiple notions of fairness like (i) *proportionality* where each of the n stakeholders receive at least $1/n$ -th of the total value of all goods (ii) *envy-freeability*, where each stakeholder weakly prefers their own bundle over anybody else’s bundle, along with other weaker notions. Variations of the problem where goods appear in an online fashion have also been studied (Aleksandrov and Walsh, 2020). A common underlying assumption in this area is that each stakeholder gets value only based on the set of items appropriated to them. In our manuscript though, with the motivation being ambulance allocation, we relax this assumption. Even if an ambulance is not placed at the base closest to a stakeholder, but say for example, at a base that is second closest to her, she could still get limited utility from the ambulance, as opposed to the case where it is placed much farther away.

Our Contributions. Our first contribution is to formally set up an abstract framework for repeatedly solving a fair allocation problem such that enhanced fairness is achieved over time, as opposed to a single round of allocation. Typically, a fair allocation could turn out to be an inefficient allocation. The framework addresses this trade-off by explicitly providing adequate control to the decision-maker on the level of the desired efficiency of the solutions.

Then, we prove that the concept of achieving fairness over time is not useful should the set of feasible allocations be a convex set, as long as the value each stakeholder obtained is a linear function of the allocation. We prove this result irrespective of the measure of fairness one uses, as long as the measure has crucial technical properties mentioned formally later. Next, we show closed-form and tight solutions for the number of rounds required for perfectly fair solutions, for special choices of the set of allocations that could be of practical interest.

However, there might also be other cases of practical interest where such closed-form solutions are harder, if not impossible, to obtain. We provide an integer programming formulation to solve these problems. The formulation hence provided can be combined directly with delayed column generation techniques in case of large instances.

With the above results and ideas, we consider the problem of allocating ambulances to multiple areas of a city, where the residents of each area are the stakeholders. This family of problems is easier than the general case as there exists an efficient greedy algorithm that provides reasonable solutions fast. However, such solutions could be far from optimal, and we provide examples where the greedy algorithm fails. Besides that, in this family of problems, we also impose restrictions on how allocations in successive rounds could be. This is reminiscent of the real-life restriction that one does not want to relocate too many ambulance vehicles each day. This added constraint makes identifying good solutions much harder as the proposed column generation-technique becomes invalid now. To solve this problem, we identify a graph induced by any candidate solution and show that deciding the feasibility of a candidate solution provided by column generation is equivalent to identifying Hamiltonian paths in this graph. Hence, we provide multiple subroutines to retain the validity of the column generation-based approach. We also provide experimental results on the computational efficiency of our method.

The paper is organized as follows. Section 2 introduces some basic definitions and concepts. Section 3 presents some theoretical proofs for simple cases, and Section 4 establishes a general framework which allows to tackle more complicated cases. This framework is used to solve a practical problem in Section 5, whose results are presented in Section 6. Finally, some conclusions are drawn in Section 7.

2 Preliminaries

In this section, we formally define the terms used throughout the manuscript. We study resource allocation in the context of fairness over time, meaning that resources are allocated to stakeholders over some time

horizon. We use \mathcal{X} to denote the set of all feasible allocations, with $n \in \mathbb{Z}_{\geq 0}$, $n \geq 2$ being the number of stakeholders among whom the allocations should be done in a fair way.

Definition 1 (Benefit function). *The benefit function $\tau : \mathcal{X} \rightarrow \mathbb{R}^n$ is a function such that the i -th component $[\tau(x)]_i$ refers to the benefit or utility enjoyed by the i -th stakeholder due to the allocation $x \in \mathcal{X}$.*

Assume ϕ to be a function that measures the unfairness associated with a given benefit pattern for the n stakeholders—a concept formalized in [Definition 2](#)—which is to be minimized. In this context, a basic fair resource allocation problem can be defined simply as

$$\min_{x \in \mathcal{X}} \phi(\tau(x)). \quad (1)$$

Many decision-makers are public servants who are either elected or nominated, and who serve in their position for a predetermined amount of time. In the context of serving the public, these officials must often juggle with limited means to ensure that everyone they serve is catered for. Public satisfaction naturally plays a role in the evaluation of their performance. Herein lies the motivation to solve resource allocation problems that consider fairness over time.

Assuming that there are T rounds of decision making, model (1) can be expanded to

$$\min_{x(t), y} \phi(y) \quad (2a)$$

$$\text{s.t. } y = \frac{1}{T} \sum_{t=1}^T \tau(x(t)) \quad (2b)$$

$$x(t) \in \mathcal{X}, \quad \forall t = 1, \dots, T. \quad (2c)$$

Here, we implicitly assume that it is sensible to add the benefits obtained at different rounds of decision making, and use the average benefit over time to compute fairness.

Now, in this section, we first define the properties that a general measure of unfairness should have for it to be considered valid in our context. To start with, we state that (i) if all the stakeholders of the problem get the same benefit or utility, then the unfairness associated with such an allocation of benefits should be 0, and (ii) the unfairness associated with benefits should be invariant to permutations of the ordering of stakeholders.

Definition 2 (Unfairness function). *Given $y \in \mathbb{R}^n$, $\phi : \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}$ determines the extent of unfairness in the allocations, if the i -th stakeholder gets a benefit of y_i . Such a function ϕ satisfies the following:*

1. $\phi(y_1, \dots, y_n) = 0 \iff y_1 = y_2 = \dots = y_n$,
2. $\phi(y_1, \dots, y_n) = \phi(\pi_1(y), \dots, \pi_n(y))$ for any permutation π .

In addition, if ϕ is a convex function, we call ϕ a convex unfairness function.

In general, trivial solutions where no benefit is obtained by any of the stakeholders, i.e., when $\tau(x)$ is a zero vector, are perfectly fair solutions! However, this results in a gross loss of efficiency.

Definition 3 (Inefficiency function). *Given \mathcal{X} and the benefit function τ , the inefficiency function $\eta : \mathcal{X} \rightarrow [0, 1]$ is defined as*

$$\eta(x) = \begin{cases} 0 & \text{if } \bar{f} = \underline{f}, \\ \frac{\bar{f} - \sum_{i=1}^n [\tau(x)]_i}{\bar{f} - \underline{f}} & \text{otherwise,} \end{cases} \quad (3)$$

where $\bar{f} = \sup_{x \in \mathcal{X}} \sum_{i=1}^n [\tau(x)]_i$, $\underline{f} = \inf_{x \in \mathcal{X}} \sum_{i=1}^n [\tau(x)]_i$, and \bar{f} and \underline{f} are assumed to be finite.

Remark 4. Note that for all feasible $x \in \mathcal{X}$, we indeed have $\eta(x) \in [0, 1]$. For the most efficient x , i.e., the x (or allocation) that maximizes the sum of benefits, $\eta(x) = 0$, while for the least efficient x , we have $\eta(x) = 1$. Thus, η serves as a method of normalization of the objective values.

We now define a single-period fair allocation problem, subject to efficiency constraints. One might always choose $\eta(x) = 1$ to retrieve the problem in model (1) without efficiency constraints.

Definition 5 (Single-period fair allocation (SPFA) problem). *Given $\bar{\eta} \in [0, 1]$, the single-period fair allocation problem is to solve*

$$\min_x \phi(\tau(x)) \quad (4a)$$

$$s.t. \quad \eta(x) \leq \bar{\eta}. \quad (4b)$$

Example 6 motivates and provides a mean of validation for the SPFA. It shows that with reasonable choices, we retrieve the intuitively fair solution.

Example 6. *Consider the case where $\mathcal{X} = \{x \in \mathbb{R}_+^n : \sum_{i=1}^n x_i = 1\}$, i.e., \mathcal{X} is a simplex. Further, assume that $[\tau(x)]_i = \tau_i x_i$ for some scalars $\tau_i > 0$, and that, w.l.o.g., $\tau_1 \geq \tau_2 \geq \dots \geq \tau_n$.*

For the choice of $\bar{\eta} = 1$, i.e., with no efficiency constraints, the solution for the SPFA problem is given by

$$x_i^* = \frac{g}{\tau_i}, \quad \forall i \quad (5a)$$

$$\text{where } g = \frac{1}{\sum_{i=1}^n \frac{1}{\tau_i}}. \quad (5b)$$

Clearly, each stakeholder enjoys a benefit of g , and hence the unfairness associated with this allocation is 0. We can notice that $\eta(x^) = \frac{\tau_1 - ng}{\tau_1 - \tau_n}$.*

Note that for the case where $n = 2$, the above simplifies to

$$\begin{aligned} x_1^* &= \frac{\tau_2}{\tau_1 + \tau_2} \\ x_2^* &= \frac{\tau_1}{\tau_1 + \tau_2} \\ \eta(x^*) &= \frac{\tau_1^2 - \tau_1 \tau_2}{\tau_1^2 - \tau_2^2}. \end{aligned}$$

We now define the fairness-over-time problem.

Definition 7 (T -period fair allocation (T -PFA) problem). *Given $\bar{\eta} \in [0, 1]$ and $T \in \mathbb{Z}_{\geq 0}$ with $T \geq 2$, the T -period fair allocation problem is to solve*

$$\min_{x(t), y} \phi(y) \quad (6a)$$

$$s.t. \quad y = \frac{1}{T} \sum_{t=1}^T \tau(x(t)) \quad (6b)$$

$$\eta(x(t)) \leq \bar{\eta}, \quad \forall t = 1, \dots, T \quad (6c)$$

$$x(t) \in \mathcal{X}, \quad \forall t = 1, \dots, T. \quad (6d)$$

We say that a fair-allocation problem (SPFA or T -PFA) has perfect fairness if the optimal objective value of the corresponding optimization problems is 0.

Example 8 (Usefulness of the T -PFA). *Let $\mathcal{X} = \{x \in \{0, 1\}^2 : x_1 + x_2 = 1\}$. Further, let $\tau(x) = (2x_1, x_2)$. Let $\bar{\eta} = 1$. Note that, in the case of the SPFA, the optimal objective is necessarily nonzero, since the benefits of the two stakeholders are unequal for every feasible solution. However, consider the 3-PFA. Now, if $x(1) = (1, 0)$, $x(2) = (0, 1)$, $x(3) = (0, 1)$, $y = (\frac{2}{3}, \frac{2}{3})$. So, for any choice of ϕ , the optimal objective value is 0, which is strictly better than the SPFA solution.*

3 Simple Cases

3.1 Convex \mathcal{X}

We have motivated in [Example 8](#) that ensuring fairness over multiple rounds could offer better results than seeking fairness in a single round. We now show that this holds only for a nonconvex \mathcal{X} . In other words, if \mathcal{X} is convex and if τ is a linear function, we necessarily get no improvement in T -PFA over SPFA.

Theorem 9. *Let f^* and f_T^* be the optimal values of SPFA and T -SPFA for some nonnegative integer T . If \mathcal{X} is convex and if τ is linear, then $f^* = f_T^*$.*

Proof. Given that τ is a linear function, we can write $[\tau(x)]_i = \tau_i^\top x$ for an appropriately chosen τ_i , for $i = 1, \dots, n$. Suppose that $x^*(t)$, for $t = 1, \dots, T$, and y^* solve the T -PFA problem. By our notation, this has an objective value of $\phi(y^*) = f_T^*$. Now, we construct a solution for the SPFA problem with an objective value equal to f_T^* . For this, consider the point $\bar{x} = \frac{1}{T} \sum_{t=1}^T x^*(t)$. First we claim that $\bar{x} \in \mathcal{X}$.

We have that $x(t) \in \hat{\mathcal{X}}$ where $\hat{\mathcal{X}} = \{x \in \mathcal{X} : \eta(x) \leq \bar{\eta}\}$. Now, we observe that $\hat{\mathcal{X}}$ is a convex set given by $\{x \in \mathcal{X} : (\sum_{i=1}^n \tau_i)^\top x \geq \bar{f} - (\bar{f} - \underline{f})\bar{\eta}\}$, where \bar{f} and \underline{f} are constants. Since \bar{x} is obtained as a convex combination of $x^*(t) \in \hat{\mathcal{X}}$, it follows that $\bar{x} \in \hat{\mathcal{X}}$. Finally, from the linearity of τ , we can set $\bar{y} = y^*$. This shows that (\bar{x}, \bar{y}) is feasible. Since $\bar{y} = y^*$, it follows that their objective function values are equal. ■ □

Having proven that multiple rounds of allocation cannot improve the fairness when \mathcal{X} is convex and τ is linear, we show that it is not necessarily due to the fact that perfect fairness is obtainable in SPFA. [Example 10](#) shows an instance where the unfairness is strictly positive with a single round of allocation, irrespective of the choice of ϕ . Naturally, due to [Theorem 9](#), perfect fairness is not possible with multiple rounds of allocation either.

Example 10. *Consider a fair allocation problem where $\mathcal{X} = \{x \in \mathbb{R}_{\geq 0}^3 : x_1 + x_2 + x_3 = 1\}$. Clearly \mathcal{X} is convex. Consider linear τ defined as*

$$\begin{aligned} [\tau(x)]_1 &= x_1 + \frac{3}{4}x_2 + \frac{3}{4}x_3 \\ [\tau(x)]_2 &= x_2 \\ [\tau(x)]_3 &= x_3. \end{aligned}$$

Let \bar{x} be a fair allocation. In such a case, we need $[\tau(\bar{x})]_1 = [\tau(\bar{x})]_2 = [\tau(\bar{x})]_3$. Thus, we need

$$\begin{aligned} g &= x_1 + \frac{3}{4}x_2 + \frac{3}{4}x_3 \\ g &= x_2 \\ g &= x_3 \end{aligned}$$

for some $g \in \mathbb{R}_+$. Solving this linear system gives the unique fair allocation as allocations of the form $(x_1, x_2, x_3) = (-0.5g, g, g)$, which necessarily violates the non-negativity constraints in the definition of \mathcal{X} . Any other allocation necessarily has $\phi(\tau(x)) > 0$.

3.2 Simplicial \mathcal{X}

We now show that perfect fairness is attainable under certain circumstances. In the theorem below, $LCM(\dots)$ refers to the least common multiple of the set of integers in its arguments.

Theorem 11. *Let $\mathcal{X} = \{x \in \mathbb{Z}_{\geq 0}^n : \sum_{i=1}^n x_i \leq a\}$ for some positive integer a . Let the benefit function be $[\tau(x)]_i = \tau_i x_i$, where each τ_i is a positive integer. Assume, w.l.o.g., that $\tau_1 \geq \tau_2 \geq \dots \geq \tau_n > 0$. Let $L = LCM(\tau_1, \dots, \tau_n)$. Then,*

1. The only perfectly fair solution within a number of periods strictly lesser than $\bar{T} = \left\lceil \frac{1}{a} \sum_{i=1}^n \frac{L}{\tau_i} \right\rceil$ is the trivial solution. i.e., $x = 0$.
2. Perfect fairness is achieved with \bar{T} periods if $\frac{\tau_1 - \min(\tau_n, \tau_1/a)}{\tau_1} \leq \bar{\eta} < 1$ and $\bar{T} > 1$.
3. Perfect fairness can be unattainable within \bar{T} periods if $\bar{\eta} < \frac{\tau_1 - \min(\tau_n, \tau_1/a)}{\tau_1}$.

Proof. Part 1. Observe that an unfairness of 0 can be achieved when for some $g \in \mathbb{R}$ we have $\sum_{t=1}^T x_i(t) = g/\tau_i$ for every $i \in \{1, \dots, n\}$. By integrality of each $x_i(t)$ and T , g must be a multiple of L , giving

$$\sum_{t=1}^T x_i(t) = \frac{\alpha L}{\tau_i}$$

for some integer $\alpha \geq 0$. Because $\bar{\eta} < 1$, α must be nonzero. Summing the above equation over $i = 1, \dots, n$, we obtain

$$\sum_{t=1}^T \sum_{i=1}^n x_i(t) = \sum_{i=1}^n \frac{\alpha L}{\tau_i}.$$

Using the fact that $\sum_{i=1}^n x_i(t) \leq a$, we obtain

$$\alpha \sum_{i=1}^n \frac{L}{\tau_i} \leq \sum_{t=1}^T a = aT \iff \frac{\alpha}{a} \sum_{i=1}^n \frac{L}{\tau_i} \leq T.$$

Since $\alpha \geq 1$, the minimum LHS is attained when $\alpha = 1$ and given that T is integer, we obtain

$$\bar{T} := \left\lceil \frac{1}{a} \sum_{i=1}^n \frac{L}{\tau_i} \right\rceil.$$

Part 2. We prove this part by exhibiting a solution feasible for \bar{T} -PFA that achieves perfect fairness over time and is hence optimal for it. Consider a reverse-lexicographic allocation method where all allocations are made to stakeholder n until the total allocation towards n sums to L/τ_n . i.e., if $L/\tau_n \leq a$, then let $x_n(1) = L/\tau_n$. Otherwise let $x_n(1) = a$ and allow $x_n(2) = \min\{a, L/\tau_n - a\}$. Repeat this process until the total allocation towards n adds up to L/τ_n . Next, repeat the same process for $n-1$ so that the total allocation towards $n-1$ adds up to L/τ_{n-1} . allocate similarly for $n-1, n-2$ and so on up to stakeholder 1. Observe that each of these allocations at any time t is in \mathcal{X} by construction. Again by construction, each player i has a value given by $\tau_i \times \frac{L}{\tau_i} = L$. If we show that the allocation in each time period has $\eta \leq \bar{\eta}$, we are done. Consider a period $\hat{t} \in \{1, \dots, \bar{T}\}$ where the allocation is the most inefficient. This would either be in the first or the last period, i.e., $\hat{t} = 1$ or $\hat{t} = \bar{T}$. This is because $T = 1$ corresponds to adding the most possible value to stakeholder n , i.e., the one who values the allocation the least. Alternatively, the most inefficient allocation could be at period T , because in the last period the constraint $\sum_{i=1}^n x_i \leq a$ might hold with a strict inequality, leading to inefficiencies. We will show in either case the inefficiency is at most $\bar{\eta}$. If the first period, η could be large as all allocations could be to n , i.e., $x(1) = (0, 0, \dots, 0, a)$. η corresponding to this allocation is $\frac{a\tau_1 - a\tau_n}{a\tau_1 - 0} = \frac{\tau_1 - \tau_n}{\tau_1} \leq \bar{\eta}$. Alternatively, in the last period, there could be minimal allocation of $\sum_{i=1}^n x_i(\bar{T}) = 1$, but in that case, it will necessarily be made to the first stakeholder, i.e., $x(T) = (1, 0, 0, \dots, 0)$. η corresponding to this allocation is $\frac{a\tau_1 - \tau_1}{a\tau_1 - 0} = \frac{\tau_1 - \tau_1/a}{\tau_1} \leq \bar{\eta}$.

Part 3. We show that perfect fairness might be unattainable with stronger efficiency requirements, by providing a family of counterexamples, each with $n = 2$ stakeholders. Consider $\hat{T} \geq 2$ and choose $\tau_1 = (\hat{T} - 1)a$ and $\tau_2 = 1$. From the first part, perfect fairness is not obtainable for $T < \bar{T} = \left\lceil \frac{1}{a} \left(\frac{(\hat{T}-1)a}{(\hat{T}-1)a} + \frac{(\hat{T}-1)a}{1} \right) \right\rceil = \hat{T}$ periods. We now show that perfect fairness is not possible with \hat{T} periods either.

Consider the set of solutions where perfect fairness is achieved in \hat{T} time steps. Any such solution should award equal total utility to both stakeholders 1 and 2. The smallest possible utility is the LCM $(\bar{T} - 1)a$, as even allocating 1 unit to player 1 already results in a utility of $\hat{T} - 1$ for them. To achieve this utility $(\hat{T} - 1)a$ allocations are required for 2. Consider the way in which such an allocation can be done. Up to permutations across periods, They are all of the form

$$\begin{aligned} x(1) &= (0, a - \delta_1) \\ x(2) &= (0, a - \delta_2) \\ &\vdots \\ x(\hat{T} - 1) &= (0, a - \delta_{\hat{T}-1}) \\ x(\hat{T}) &= \left(\sum_{i=1}^{\hat{T}-1} \delta_i, 1 \right) \end{aligned}$$

for some integers $\delta_i \in \{1, 2, \dots, a\}$ satisfying $0 \leq \sum_{i=1}^{\hat{T}-1} \delta_i \leq a - 1$. From the hypothesis, we need $\bar{\eta} < \frac{\tau_1 - \min(\tau_n, \tau_1/a)}{\tau_1} = \frac{(\hat{T}-1)a-1}{(\bar{T}-1)a}$.

Now consider the inefficiency function of the allocation $x(1)$. This is minimized when $\delta_1 = 0$. In that case, $\eta(x(1)) = \frac{(\hat{T}-1)a-1+\delta_1/a}{\hat{T}-1} = \frac{(\hat{T}-1)a-1}{(\bar{T}-1)a} > \bar{\eta}$ shows that it is infeasible. Thus, any perfectly fair allocation for \bar{T} -PFA problem is an infeasible allocation, providing the necessary counterexample. ■ □

The first part of [Theorem 11](#) states that for any n and a \mathcal{X} of the specified form, perfect fairness is possible after a specific, finite number of periods, \bar{T} , provided the efficiency requirements are as stated. The second part shows tightness of the results saying, for any other stricter efficiency requirements, there exists a counterexample with just two stakeholders, such that perfect fairness is impossible in \bar{T} steps.

Example 12. Consider [Theorem 11](#) applied to [Example 8](#). We have $\mathcal{X} = \{(1, 0), (0, 1)\}$, and thus $a = 1$. Then, $\tau(x) = (2x_1, x_2)$, so $\tau_1 = 2$ and $\tau_2 = 1$, and $\tau_1 \geq \tau_2 \geq 0$ holds. In addition, $L = \text{LCM}(2, 1) = 2$. Then, $\bar{T} = \left\lceil \frac{1}{a} \sum_{i=1}^n \frac{L}{\tau_i} \right\rceil = \left\lceil \frac{1}{1} \left(\frac{2}{2} + \frac{2}{1} \right) \right\rceil = 3$.

3.3 Sparse Simplicial \mathcal{X}

The sparse simplicial form for \mathcal{X} is another interesting case from a practitioner's perspective. For example, a government might want to allot money to n possible projects, but if it divides among all of them, then it might be insufficient for any of them. So, there could be a restriction that the government allots it to at most r projects.

Theorem 13. Let $\mathcal{X} = \{x \in \mathbb{R}_{\geq 0}^n : \sum_{i=1}^n x_i \leq a, \|x\|_0 \leq r\}$ where $\|\cdot\|_0$ is the sparsity pseudo-norm, which counts the number of non-zero entries in its argument. Assume that the benefit function is $[\tau(x)]_i = \tau_i x_i$, where each τ_i is a positive integer. Assume, w.l.o.g., that $\tau_1 \geq \tau_2 \geq \dots \geq \tau_n > 0$. We denote by (p, q) the quotient and the remainder for $\frac{n}{r}$. Let $\bar{T} = \left\lceil \frac{n}{r} \right\rceil$ and $1 > \bar{\eta} \geq \frac{a\tau_1 - qv}{a\tau_1}$ where $v = a / \max \left\{ \sum_{i=(\bar{T}-1)r+1}^n \frac{1}{\tau_i}, \sum_{i=(p-1)r+1}^{pr} \frac{1}{\tau_i} \right\}$. There is no perfect assignment if $T < \bar{T}$ whereas perfect fairness is attainable in \bar{T} periods.

Proof. First observe that at most r stakeholders can be provided value in any given period. Thus except the case where everybody gets a value of 0, perfect fairness is impossible in any fewer than $\bar{T} = \left\lceil \frac{n}{r} \right\rceil$ periods.

Now we show that perfect fairness is possible in \bar{T} periods, if the conditions in the theorem statement are satisfied. Consider the following allocation: In period t for $1 \leq t \leq \bar{T} - 1$, allocate $\frac{v}{\tau_i}$ for $i = (t-1)r + 1, (t-1)r + 2, \dots, tr$ and 0 to the rest. For $t = \bar{T}$, allocate $\frac{v}{\tau_i}$ for $i = (\bar{T}-1)r + 1, \dots, n$. Note that in each

period $t < \bar{T}$, we allocate exactly to r stakeholders, and for $t = \bar{T}$, if $q = 0$, we allocate to r stakeholders, else we allocate to q stakeholders, thus always satisfying $\|x\|_0 \leq r$.

In any period $t < \bar{T}$, note that $\sum_{i=1}^n x_i(t) = \sum_{i=(t-1)r+1}^{tr} \frac{v}{\tau_i} \leq v \sum_{i=(p-1)r+1}^{pr} \frac{1}{\tau_i} \leq v \max \left\{ \sum_{i=(\bar{T}-1)r+1}^n \frac{1}{\tau_i}, \sum_{i=(p-1)r+1}^{pr} \frac{1}{\tau_i} \right\}$ a satisfying the inequality constraint. For $t = \bar{T}$, we have $\sum_{i=1}^n x_i(t) = \sum_{i=(\bar{T}-1)r+1}^n \frac{v}{\tau_i} \leq v \max \left\{ \sum_{i=(\bar{T}-1)r+1}^n \frac{1}{\tau_i}, \sum_{i=(p-1)r+1}^{pr} \frac{1}{\tau_i} \right\}$ a, again satisfying the inequality constraint.

We can ensure that in the first $\bar{T} - 1$ periods, each of the allocated players receives a value of v , and hence the total benefit of the allocation is rv , giving the $\eta(x(t)) = \frac{a\tau_1 - rv}{a\tau_1} \leq \bar{\eta}$ since $q < r$. In the last period, we allocate a utility of v to q players if $q > 0$ else to r players. The total benefit obtained in this round is at least qv . Feasibility follows as before and $\eta(x(\bar{T})) = \frac{a\tau_1 - qv}{a\tau_1} \leq \bar{\eta}$, which is feasible. Perfect fairness follows since the benefit to each player is v . ■. □

Remark 14. Unlike the setting in [Theorem 11](#), [Theorem 13](#) is not tight with respect to $\bar{\eta}$. In other words, it remains unknown whether decreasing the allowed value of inefficiency $\bar{\eta}$ still allows one to achieve perfect fairness in \bar{T} periods.

Remark 15. We note that the above results, [Theorems 9, 11 and 13](#), are agnostic to the choice of ϕ , and only use the fact that $\phi(x_1, \dots, x_m) = 0 \iff x_1 = \dots = x_n$. Any result that holds for $\phi(\cdot) \neq 0$ has to necessarily depend upon the choice of ϕ .

Example 16. Consider [Theorem 13](#) applied to the following setting. $\mathcal{X} = \{x \in \mathbb{R}_{\geq 0}^3 : x_1 + x_2 + x_3 \leq 1, \|x\|_0 \leq 2\}$. Let $\tau_1 = 5$ and $\tau_2 = 3$. Now $\bar{T} = \lceil \frac{n}{r} \rceil = \lceil \frac{3}{2} \rceil = 2$. One can calculate that $v = 2\frac{2}{9}$, $\bar{\eta} = \frac{5}{9}$. The allocation $(\frac{4}{9}, \frac{5}{9}, 0)$ in period 1 and $(0, 0, \frac{740}{999})$ in period 2, gives all players a utility of $2\frac{2}{9}$ achieving perfect fairness.

4 General Combinatorial \mathcal{X}

In many practical areas of interest, \mathcal{X} could be more complicated than the sets presented in [Section 3](#). Let $\mathcal{X} = \{x^1, x^2, \dots, x^k\}$. We assume that \mathcal{X} is finite, but typically has a large number of elements, given in the form of a solution to a combinatorial problem. We consider a benefit function where $\tau(x) = \Gamma x$, for some matrix Γ . We thus have $[\tau(x)]_i = \tau_i^\top x$ for $i = 1, \dots, n$. The efficiency constraint here is trivial, as the inefficient allocations x^j can be removed from \mathcal{X} to retain another (smaller) finite \mathcal{X} . In fact, with linear τ , the efficiency constraint is a linear inequality of the form $\sum_{i=1}^n \tau_i^\top x \geq \bar{f} - \bar{\eta}(\bar{f} - \underline{f})$.

Let q_j be the number of times the allocation $x^j \in \mathcal{X}$ is chosen over T rounds of decision. In this setting, the T -PFA problem can be restated as¹

$$\min_{q, y} \quad \phi(y) \tag{7a}$$

$$\text{s.t.} \quad y_i = \frac{1}{T} \tau_i^\top \left(\sum_{j=1}^k q_j x^j \right), \quad \forall i = 1, \dots, n \tag{7b}$$

$$\sum_{j=1}^k q_j = T \tag{7c}$$

$$q_j \in \mathbb{Z}_{\geq 0}, \quad \forall j = 1, \dots, k. \tag{7d}$$

Note that since the theorems of [Section 3](#) do not extend to this case, and since [Example 10](#) shows that a perfectly fair solution may not even exist, we know neither the value of the fairest feasible solution, nor the smallest value of T that guarantees such a solution. We now present a two-phase integer program that finds

¹The reader may have observed that an approach based on column generation would lend itself well for such a model. This is discussed in [Section 5](#).

the smallest value of T that guarantees the fairest feasible solution. In the first phase, we are interested in finding the fairest feasible solution. This is achieved by solving model (7), replacing (7c) with $\sum_{j=1}^k q_j \geq 1$ to ensure that not only the fairest solution is returned, but that this solution not be $(0, 0, \dots, 0)$. Let ϕ^* be the value of the solution found in the first phase. In the second phase, we are interested in finding the smallest T which can accommodate a solution of value ϕ^* . This is achieved by solving model (7), again replacing (7c) with $\sum_{j=1}^k q_j \geq 1$, adding $\phi(y) = \phi^*$, and changing the objective to $\min_{q,y} \sum_{j=1}^k q_j$.

We should note that the value of T^* found by the second phase could be quite high, even for simple choices of \mathcal{X} . If perfect fairness can be attained in a time horizon of size, say, 10000, and that a discrete time point represents a day, then this period of 27-odd years would probably not be suitable for most practical applications, since it is only by reaching the end of the time horizon that optimal fairness is guaranteed. In practice, then, large time horizons can be problematic. This motivates the need for consistent fairness on a smaller scale.

Some compromises can be made which would mitigate this issue. Manually fixing T to a value not higher than the expected duration of the problem would ensure a fair distribution of resources, since reaching the end of the time horizon would be assured. Accepting some degree of unfairness would also decrease the length of the time horizon and achieve similar results. Example 17 illustrates these two options.

Example 17. Consider a fair allocation problem where $\mathcal{X} = \{x \in \mathbb{Z}_{\geq 0}^2 : x_1 + x_2 = 1\}$. Further consider τ defined as

$$\begin{aligned} [\tau(x)]_1 &= x_1 + \frac{15}{37}x_2 \\ [\tau(x)]_2 &= x_2 + \frac{15}{47}x_1. \end{aligned}$$

Perfect fairness can be achieved when $T = 1109$, but if we fix $T = 5$, the difference between the benefits of both stakeholders will be $\sim 13\%$, and if we allow a small difference of 0.1% between their benefits, this could be achieved when $T = 15$.

Another compromise would be to choose a fair way of reaching the solution. Given that we know the value of T^* which grants optimal fairness, we could then identify a *best path* leading to the optimal solution. A best path is the one which ensures that unfairness is kept as low as possible in the intermediate time points, without sacrificing optimality at the end of the time horizon. This would, however, require solving another linear program, which may prove burdensome if the time horizon were too large.

One may think that greedily aiming for the fairest solution at each time point—while considering any unfairness incurred in previous time points, and enforcing minimal efficiency constraints by ensuring that all available resources are used at each time point—could lead to optimality at the end of the time horizon. Theorem 18 shows that it is not the case.

Theorem 18. A greedy approach does not provide any guarantee of optimality for the T -PFA problem.

Proof. Consider a fair allocation problem where $\mathcal{X} = \{x \in \mathbb{Z}_{\geq 0}^3 : x_1 + x_2 + x_3 = 1\}$. Further consider τ defined as

$$\begin{aligned} [\tau(x)]_1 &= \epsilon x_1 + \frac{1}{2}x_2 + \frac{1}{2}x_3 \\ [\tau(x)]_2 &= x_2 \\ [\tau(x)]_3 &= x_3. \end{aligned}$$

where $0 < \epsilon < \frac{1}{2}$. If $T = 2$, the only perfectly fair solution (and, as it happens, the most efficient too) is achieved by allocating the resource in turn to x_2 and x_3 , over T . Yet, for any reasonable measure of unfairness, the greedy choice would be to allocate the resource twice to x_1 . For instance, suppose that ϕ was defined by the difference between the largest and smallest benefits, i.e., $\phi(x) = \max_{i \in \{1, \dots, n\}} [\tau(x)]_i - \min_{i \in \{1, \dots, n\}} [\tau(x)]_i$.

In such a case, the greedy choice could never deviate from a continuous allocation to x_1 , as any other allocation would always be less fair—here, the solution would always be $T\epsilon - 0 = T\epsilon$. Thus, a greedy approach does not necessarily lead to an optimal solution. ■ □

We should note that there is a trade-off between fairness and efficiency: As illustrated in [Example 19](#), enforcing constraints to increase fairness will generally decrease efficiency, and vice-versa.

Example 19. Consider a SPFA problem where $\mathcal{X} = \{x \in \mathbb{Z}_{\geq 0}^3 : 1 \leq x_1 + x_2 + x_3 \leq 3\}$. Further consider τ defined as

$$\begin{aligned} [\tau(x)]_1 &= x_1 + \frac{1}{2}x_2 + \frac{1}{2}x_3 \\ [\tau(x)]_2 &= x_2 + \frac{1}{2}x_1 \\ [\tau(x)]_3 &= x_3 + \frac{1}{2}x_1. \end{aligned}$$

By enforcing $\eta(x) = 0$, the only feasible solution is $x = (3, 0, 0)$ which gives $\tau(x) = (3, \frac{3}{2}, \frac{3}{2})$ which is not perfectly fair. In contrast, by necessitating perfect fairness, the only feasible solution is $x = (0, 1, 1)$ which gives $\tau(x) = (1, 1, 1)$. This corresponds to $\eta(x) = \frac{6-3}{6-(3/2)} = \frac{2}{3}$.

We should also note that the value of T may have a noticeable impact on both fairness and efficiency, as illustrated in [Example 20](#).

Example 20. Consider a fair allocation problem where $\mathcal{X} = \{x \in \mathbb{Z}_{\geq 0}^4 : 1 \leq x_1 + x_2 + x_3 + x_4 \leq 3\}$. Further consider τ defined as

$$\begin{aligned} [\tau(x)]_1 &= x_1 + x_2 \\ [\tau(x)]_2 &= x_2 + x_1 \\ [\tau(x)]_3 &= x_3 + x_4 \\ [\tau(x)]_4 &= x_4 + x_3. \end{aligned}$$

The smallest T accommodating perfect fairness is $T = 1$ with $x = (1, 0, 1, 0)$ ² for $\tau(x) = (1, 1, 1, 1)$ and an inefficiency of $\eta(x) = \frac{6-4}{6-2} = 0.5$ (with one wasted resource). However, $T = 2$ also accommodates solutions with perfect fairness, albeit with an increased efficiency. Take, for instance, $x(1) = (3, 0, 0, 0)$ and $x(2) = (0, 0, 0, 3)$ over two time points, for a combined $\tau(x(1)) + \tau(x(2)) = (3, 3, 3, 3)$ and an inefficiency of $\eta(x(t)) = 0$ in each period t . This represents a noticeable increase in efficiency over choosing the initial solution twice to cover the same time horizon.

In other words, reaching perfect fairness in the shortest possible horizon might not lead to the most efficient solutions.

5 Ambulance Location and Relocation

The task of allocating ambulances to a set of bases in a region is a well-known problem in operations research ([Brotcorne et al., 2003](#)). The objective of this problem is generally one of efficiency: Ambulances should be allocated to prime spots such that they can quickly reach the maximum number of people.

The definition of this problem is not unified across the literature ([Brotcorne et al., 2003](#), [Gendreau et al., 2001](#)).³ Most definitions, however, agree that the population should receive an efficient service, which is generally translated into some form of coverage. In this paper, we are not solely concerned by efficiency, but we are further interested in fairness: Ambulances should be allocated such that the same set of people is not always at a disadvantage with respect to access to a quick service.

²Or any other equivalent solution.

³This is due to the fact that research projects often work with datasets associated with specific regions, each of which must follow local guidelines and regulations.

5.1 Preliminary Problem Formulation

In this manuscript, we adopt the following definition. The region to which ambulances are allocated is divided into n demand zones, the travel time between each pair being a known quantity. Ambulances can only be allocated to bases, which are located within a subset of the zones. Each zone (equivalently, the people living in each zone) are individual stakeholders in this model. A pre-chosen T rounds of decision is modeled to occur, over which the ambulances should be allocated in a fair and efficient manner. In each round, a configuration x of how m ambulances are placed is chosen. Next, we define the benefit function for the zones. In this model, each zone i , based on its population, has a demand of ζ_i ambulance vehicles. A zone i is said to be *covered* in configuration x if there are at least ζ_i ambulances located in bases from which zone i could be reached in a time less than a chosen time limit called the *response threshold time*. Now, $[\tau(x)]_i = 1$ if zone i is covered by the configuration x and $[\tau(x)]_i = 0$ otherwise. We note that we use a nonlinear benefit function in this case, as opposed to the simpler cases in [Section 3](#) to both reflect the fact that a larger number of ambulances are essential to sufficiently serve certain regions and also to demonstrate the robustness of our methods to even certain classes of nonlinear (piecewise linear) benefit functions. The nonlinearity, as shown below, can be modeled using auxiliary variables and linear functions, to conform to the form of [\(7\)](#). Furthermore, efficiency constraints analogous to [\(6c\)](#) are enforced by stating that at least a fraction f of the zones should be covered in each allocation. The unfairness metric ϕ used is the difference between the two zones which are most often and least often covered, over T , i.e., $\phi(\tau) := \min_{g,h} \{g - h : g \geq \tau_i \geq h, \forall i \in \{1, \dots, n\}\}$.

With the above, the problem is a standard T-PFA problem in [Definition 7](#). However, the ambulance allocation problem involves additional constraints on allocations between two consecutive periods. Decision-makers prefer policies that ensure that not too many ambulances have to shift bases on a daily basis. Thus, between two consecutive allocations, we would ideally like to have not more than a fixed number r of ambulances to shift bases. We refer to this constraint as the *transition constraint*.

Variables	
$\phi(y)$	Unfairness associated with benefits y
y_i	Average benefit of zone i
$x_i(t)$	Number of ambulances at base i on time t
$v_i(t)$	Number of ambulances that can reach zone i at time t
$\tau_i(t)$	Benefit of zone i at time t
Parameters	
\mathcal{B}	Set of bases
\mathcal{T}	Set of time points
n	Number of zones
m	Number of ambulances
a_{ji}	1 if zones i and j are connected, 0 otherwise
ζ_i	Demand of ambulances for zone i
T	Time horizon
f	Fraction of zones that must be covered
r	Number of ambulances allowed to shift bases

Table 1: Variables and parameters of the AWT.

Now, the problem can be formally cast as follows:

$$\min_{y, v, x, \tau, \phi} \quad \phi(y) \tag{8a}$$

$$\text{s.t.} \quad x_i(t) = 0 \quad \forall i \notin \mathcal{B}; \forall t \in \mathcal{T} \tag{8b}$$

$$\sum_{i=1}^n x_i(t) \leq m \quad \forall t \in \mathcal{T} \tag{8c}$$

$$v_i(t) = \sum_{j=1}^n a_{ji} x_j(t) \quad \forall i = 1, \dots, n; t \in \mathcal{T} \quad (8d)$$

$$v_i(t) \leq (\zeta_i - 1) + m \tau_i(t) \quad \forall i = 1, \dots, n; t \in \mathcal{T} \quad (8e)$$

$$v_i(t) \geq \zeta_i - m(1 - \tau_i(t)) \quad \forall i = 1, \dots, n; t \in \mathcal{T} \quad (8f)$$

$$y_i = \frac{1}{T} \sum_{t=1}^T \tau_i(t) \quad \forall i = 1, \dots, n \quad (8g)$$

$$\sum_{i=1}^n \tau_i(t) \geq fn \quad \forall t \in \mathcal{T} \quad (8h)$$

$$x_i(t) \in \mathbb{Z}_{\geq 0} \quad \forall i = 1, \dots, n; t \in \mathcal{T} \quad (8i)$$

$$\tau_i(t) \in \{0, 1\} \quad \forall i = 1, \dots, n; t \in \mathcal{T} \quad (8j)$$

$$\sum_{i=1}^n |x_i(t+1) - x_i(t)| \leq 2r \forall t \in 1, \dots, T-1 \quad (8k)$$

Here, $\mathcal{T} = \{1, \dots, T\}$, $x_i(t)$ is the number of ambulances allotted to zone i at time t . Constraints (8b) ensure that allocation occurs only if i is an ambulance base. Here, \mathcal{B} is the set of ambulance bases. Constraints (8c) limit the number of ambulances available for allocation in each round. Moreover, a_{ij} is a binary parameter which is 1 if an ambulance can go from i to j within the response threshold time and 0 otherwise. Through constraints (8d), $v_i(t)$ counts the number of ambulances that can reach zone i within the response threshold time. Here, $\tau_i(t)$ is 1 if $v_i(t)$ is at least ζ_i , i.e., if the demand of ambulances in zone i is satisfied, and 0 otherwise. This is accomplished by constraints (8e) and (8f). Constraints (8h) take care of efficiency and ensure only allocations that cover at least a fraction f of all zones are to be considered. Finally, constraints (8k) are the transition constraints, which can easily be reformulated through linear inequalities. It ensures that not too many ambulances shift bases between consecutive time periods.

The problem in (8) is a mixed-integer linear program (MILP). However, the problem is symmetric with respect to certain permutations of the variables. Symmetry makes it particularly hard for modern branch-and-bound-based solvers (Margot, 2010). Even if one relaxes the transition constraints (8k), the symmetry exists. One can observe that, relaxing the transition constraints (8k) in (8), we have a T-PFA problem. We call this relaxed problem defined by (8a) to (8j) as the Ambulance-without-transition constraints (AWT) problem.

5.2 A Branch-and-Price Reformulation of the AWT

Since the AWT in (8a) to (8j) is a T-PFA problem, it can be readily written in the form shown in (7) and hence can be solved using branch-and-price. First, we note that without the constraint in (8k), any feasible solution or optimal solution remains feasible or optimal after permutations to t . Thus, the following version of the problem could be used to solve the relaxed problem without the symmetry. Here, we only *count* the number of times each configuration might be used over T periods.

Variables			Parameters	
g	Largest average benefit	τ_{ij}	Benefit of zone i in configuration j	
h	Smallest average benefit	T	Time horizon	
y_i	Average benefit of zone i	n	Number of zones	
q_j	Number of times configuration j is used	k	Number of configurations	

Table 2: Variables and parameters of the MP.

The Master Problem (MP).

$$\min \quad g - h \quad (9a)$$

$$\text{s.t.} \quad g \geq y_i \quad (\alpha_i) \quad \forall i = 1, \dots, n \quad (9b)$$

$$y_i \geq h \quad (\beta_i) \quad \forall i = 1, \dots, n \quad (9c)$$

$$y_i = \frac{1}{T} \sum_{j=1}^k \tau_{ij} q_j \quad (\lambda_i) \quad \forall i = 1, \dots, n \quad (9d)$$

$$\sum_{j=1}^k q_j = T \quad (\mu) \quad (9e)$$

$$q_j \geq 0 \quad \forall j = 1, \dots, k \quad (9f)$$

$$q_j \in \mathbb{Z} \quad \forall j = 1, \dots, k. \quad (9g)$$

where q_j counts the number of times the configuration defined by x^j is used. The benefit obtained by stakeholder i due to the configuration x^j is τ_{ij} . y_i is the average benefit that stakeholder i enjoys through the time horizon of planning. The objective (9a) is to minimize the difference between the largest (g) and the smallest (h) average benefits. Note that, once we solve the MP, an equivalent solution to the AWT could be obtained by arbitrarily considering the allocations x^j for q_j times in $x(1), \dots, x(T)$ of the AWT. Similarly, given a solution to the AWT, one could immediately identify a corresponding solution to the MP.

However, since the number of configurations are typically exponentially large in n and m , and we only might use a handful of them in a solution, we could resort to a branch-and-price approach where the MP only contains a subset of the configurations.

Referring to the continuous relaxation of MP as CMP, the dual of CMP can be found in Appendix A.

Considering only a subset of columns in the CMP is equivalent to considering only a subset of constraints in (13e). Given some dual optimal solution $(\alpha^*, \beta^*, \lambda^*, \mu^*)$ to the dual of the restricted CMP, one can find the most violated constraint in (13e) and include the corresponding column in the restricted CMP.

The Pricing Problem.

$$\min \quad \sum_{i=1}^n \tau_i \lambda_i^* \quad (10a)$$

$$\text{s.t.} \quad (\tau_i = 1) \iff (a_i^T x \geq \zeta_i), \quad \forall i = 1, \dots, n \quad (10b)$$

$$(\tau_i = 0) \iff (a_i^T x \leq \zeta_i - 1), \quad \forall i = 1, \dots, n \quad (10c)$$

$$\sum_{i=1}^n x_i \leq m, \quad \forall i \in \mathcal{B} \quad (10d)$$

$$\sum_{i=1}^n \tau_i \geq fn \quad (10e)$$

$$\tau \in \{0, 1\}^n \quad (10f)$$

$$x \in \mathbb{Z}_{\geq 0}^n. \quad (10g)$$

The minimal efficiency constraints are embedded in the pricing problem. The time horizon is fixed in (9e).

Constraints (10b) and (10c) help compute the benefits to each zone, τ and can clearly be rewritten with integer variables and linear constraints. No more than m ambulances may be used (10d), and the configuration must cover at least a fraction f of the zones (10e).

The MP reformulation of the AWT in (8a) - (8j), can now be solved very efficiently using branch-and-price. However, a solution thus obtained might violate (8k). Further, given the solution in the space of variables in the MP it is not immediate if one can easily verify whether the constraint (8k) is or is not satisfied. So, given a feasible point for the MP, we define the configuration graph as follows, and then show that the point satisfies (8k) if and only if the configuration graph has a Hamiltonian path.

5.3 Checking (8k)

In the previous section, we proposed a branch-and-price to solve the relaxed version of the AWT without including constraints (8k). First, we provide an algorithm to check the feasibility of a solution provided by branch-and-price in Theorem 22. If it is feasible, then we are trivially done. If the solution is infeasible, we provide routines to “cut-off” such an infeasible solution and continue the branch-and-price algorithm. Meanwhile, we also show how using cutting planes could be impractical in the space of solutions considered in the MP or a binarized version of the problem. Then, we talk about a three-way branching scheme that could work. However, we finally resort to constraint programming, due to the limitations of commercial solvers in implementing three-way branching.

Definition 21 (Configuration graph). *Given a solution \bar{q} to the MP, define $\mathcal{I} = \{j : \bar{q}_j \geq 1\}$. The configuration graph (CG) is defined as $G = (V, E)$, where $V = \{(j, j') : j \in \mathcal{I}; j' \in \{1, 2, \dots, \bar{q}_j\}\}$ and $E = \{((j_1, j'_1), (j_2, j'_2)) : \|\bar{x}^{j_1} - \bar{x}^{j_2}\|_1 \leq 2r\}$ where \bar{x}^{j_1} and \bar{x}^{j_2} are the configurations corresponding to the variables \bar{q}_{j_1} and \bar{q}_{j_2} .*

Theorem 22. *Given a point \bar{q} that is feasible to the MP, there exists a corresponding \bar{x} that is feasible to the AWT if the CG defined by \bar{q} contains a Hamiltonian path. Conversely, if there is a feasible solution to the AWT which uses only the configurations with indices in $\mathcal{I} = \{j : \bar{q}_j \geq 1\}$, then the corresponding CG has a Hamiltonian path.*

Proof. Observe that G has exactly T vertices, since if a configuration is found more than once in \bar{q} , it is split into as many distinct vertices in V , which are distinguished by the second element of the tuple. An edge $((u_1, u_2), (v_1, v_2))$ in E indicates that movement between configurations indexed by u_1 and v_1 does not violate the transition constraints. A Hamiltonian path is a path that visits each vertex exactly once. As such, any Hamiltonian path in G proves that a feasible sequence of transitions which does not violate the transition constraints exists between the configurations in \bar{q} . Given a Hamiltonian path $(v_1, w_1), (v_2, w_2), \dots, (v_T, w_T)$ in G , a feasible sequence $x(1), x(2), \dots, x(T)$ for the AWT would be $\bar{x}^{v_1}, \bar{x}^{v_2}, \dots, \bar{x}^{v_T}$.

Conversely, given a feasible sequence $x^{j_1}, x^{j_2}, \dots, x^{j_T}$ for the AWT, a Hamiltonian path can be constructed for G as follows. The t -th vertex of the path is (j_t, β) where $\beta = 1 +$ the number of times the configuration denoted by x^{j_t} has appeared in the first $t - 1$ vertices of the path. ■ □

Following Theorem 22, one can construct the CG with exactly T vertices and can check if the solution to the MP is feasible to the AWT. If yes, we are done. If not, the way we can proceed is detailed in the rest of this section.

5.3.1 Cutting Planes and Binarization

The most natural way to eliminate a solution that does not satisfy a constraint is by adding a cutting plane. This is a common practice in the MILP literature.

In cases where a more sophisticated cutting plane is not available, but every feasible solution is determined by a binary vector, no-good cuts could be used to eliminate infeasible solutions one by one (D’Ambrosio et al., 2010). However, in our problem of interest, the variables x in the MP are general integer variables. It is possible that a point x that violates the transition constraint could lie strictly in the convex hull of solutions that satisfy the transition constraint. Hence it could be impossible to separate x using a cutting plane. Example 23 demonstrates the above phenomenon.

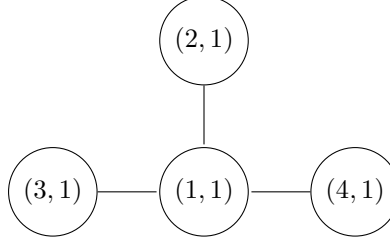


Figure 1: Configuration graph for the solution in [Example 23](#).

Example 23. Consider the case where the MP has an optimal solution given by $q = (q_1, q_2, q_3, q_4) = (1, 1, 1, 1)$ and the allocations x_1, x_2, x_3, x_4 (in the context of the AWT) corresponding to q_1, q_2, q_3, q_4 are $(3, 3, 3, 3)$, $(4, 2, 3, 3)$, $(2, 4, 3, 3)$, $(3, 3, 2, 4)$ respectively. Let $r = 1$. Then, the CG corresponding to this solution is a tree as shown in [Figure 1](#) and hence does not have a Hamiltonian path. On the other hand, for the choices q^i, q^{ii}, q^{iii} and q^{iv} being $(4, 0, 0, 0)$, $(0, 4, 0, 0)$, $(0, 0, 4, 0)$, $(0, 0, 0, 4)$, the CGs are all K_4 , i.e., complete graphs and hence have a Hamiltonian path, trivially.

Now, it is easy to see that q lies in the convex hull of q^i, q^{ii}, q^{iii} and q^{iv} and while each of the latter solutions satisfies the transition constraint, q does not. Thus no valid cutting plane can cut only the infeasible point.

While [Example 23](#) shows that an infeasible solution cannot always be separated using cutting planes, it could still be possible that there is an extended formulation where the infeasible point could be separated.

Naive Binarization. A natural choice of an extended formulation comes from binarizing each integer variable q_j in the MP. This is possible because each q_j is bounded above by T and below by 0. Thus one can write constraints of the form $q_j = b_j^1 + 2b_j^2 + 4b_j^4 + 8b_j^8 + \dots$ where the summation extends up to the smallest ℓ such that $2^\ell > T$, i.e., $\ell = \lceil \log_2 T \rceil$. If each b_j^ℓ is binary, any integer between 0 and T can be represented as above. Having written the above binarization scheme, one can separate any solution q by adding a no-good cut on the corresponding binary variables. A no-good cut is a linear inequality that separates a single vertex of the 0 – 1 hypercube without cutting off the rest ([Balas and Jeroslow, 1972](#), [D’Ambrosio et al., 2010](#)). An example is shown in [Example 24](#).

Example 24. Consider the problem in [Example 23](#). Binarization to separate the solution $q = (1, 1, 1, 1)$ can be done by adding the following constraints to the MP.

$$q_j = b_j^1 + 2b_j^2 + 4b_j^4 \quad \forall j = 1, \dots, 4 \quad (11a)$$

$$\sum_{j=1}^4 (b_j^1 + (1 - b_j^2) + (1 - b_j^4)) \leq 3 \quad (11b)$$

Note that the second constraint above (the no-good constraint) is violated only by the binarization corresponding to the solution $u = (1, 1, 1, 1)$ and no other feasible solution is cut off.

The potential downside with the above scheme is that one might have to cut off a prohibitively large number of solutions before reaching the optimal solution. And with the column generation introducing new q -variables, this could lead to an explosion in the number of new variables as well as the number of new constraints.

Strengthened Binarization. When the CG corresponding to q is not just lacking a Hamiltonian path, but is a disconnected graph (a stronger property holds), one could add a stronger cut, which could potentially cut off multiple infeasible solutions. In this procedure, we add a binary variable b_j for each q_j such that $b_j = 1$ if and only if $q_j \geq 1$. Now, observe that if a CG corresponding to the solution \bar{q} is disconnected, then

it will be disconnected for all q whose non-zero components coincide with the non-zero components of \bar{q} . i.e., no solution with the same support as that of \bar{q} could satisfy the transition constraints. Thus one could add a no-good cut on these b_j binary variables, which cuts off all the solutions with the same support as \bar{q} .

Unlike naive binarization, while this cuts off multiple solutions simultaneously, it could happen that the CG is connected, but just does not have a Hamiltonian path. In such a case, no cut could be added by the strengthened binarization scheme, and one might have to resort to the naive version.

5.3.2 Three-Way Branching

An alternative to the naive binarization scheme is three-way branching. While this could work as a stand-alone method, it could also go hand in hand with the strengthened binarization mentioned earlier. This method takes advantage of the three-way branching feature that some solvers, for example, SCIP (Gamrath et al., 2020a,b), have.

In this method, as soon as a solution \bar{q} satisfying all the integrality constraints but violating the transition constraint is found, the following actions are performed. First, if the lower bound and the upper bound for every component of \bar{q} match, then we are at a leaf that can be discarded as infeasible. If not, find a component j (in our case, a configuration j) such that \bar{q}_j is strictly different from at least one of the bounds. Now, we do a three-way branching on the variable q_j where the new constraints in each of the three branches are (i) $q_j \leq \bar{q}_j - 1$ (ii) $q_j = \bar{q}_j$ (iii) $q_j \geq \bar{q}_j + 1$.

Finite termination follows from the fact that \bar{q}_j is infeasible for branches (i) and (iii) and hence will never be visited again. For branch (ii), we have q_j such that its lower and upper bounds are equal to \bar{q}_j . Thus, we have one more fixed variable and this variable will never be branched on again.

Three-way branching is used as opposed to regular two-way branching but on integer variables because of the following reasons. First, branching with (i) $q_j \leq \bar{q}_j - 1$ (ii) $q_j \geq \bar{q}_j + 1$ is invalid, as it could potentially cut off other feasible solutions with $q_j = \bar{q}_j$ but the solution differing in components other than j . Branching with (i) $q_j \leq \bar{q}_j - 1$ (ii) $q_j \geq \bar{q}_j$ could cycle, as we have not fixed any additional variable in the second branch, nor have we eliminated the infeasible solution. Thus, the LP optimum in the second branch is going to be \bar{q} again, and cycling ensues.

5.3.3 Constraint Programming

The final alternative we can use to enforce transition constraints (8k) is constraint programming (CP). Namely, CP is a programming paradigm for solving combinatorial problems. A CP model is defined by a set of variables, each of which is allowed values from a finite set, its *domain*. The relationship between these variables is determined by the constraints of the problem. These succinct constraints can encapsulate complex combinatorial structures, such as the packing of items into bins, for example. A solver then solves the problem by enforcing consistency between the variables and the constraints, and using branching and backtracking techniques to explore the solution space. Before defining the CP model to enforce constraints (8k), we define the compact configuration graph, and explain its relationship with the CG.

Definition 25 (Compact configuration graph). *Given a feasible solution q^* to the continuous relaxation of the MP (CMP), the compact configuration graph (CCG) is defined as $G = (V, E)$, where $V = \mathcal{A} := \{j : q_j^* > 0\}$ and $E = \{(v, w) : \|\bar{x}^v - \bar{x}^w\|_1 \leq 2r, \forall v, w \in \mathcal{A}\}$.*

Definition 26 (Walk). *A walk in an undirected graph $G = (V, E)$ is a finite sequence of vertices v_1, \dots, v_k , not necessarily distinct, such that for each i , $(v_i, v_{i+1}) \in E$.*

Theorem 27. *Every walk of length T in a CCG constructed from a solution q^* to CMP corresponds to a feasible solution of the AWT.*

Proof. Let $G = (V, E)$ be the CCG given a solution q^* . Let $W = v_1, v_2, \dots, v_T$ be a walk of length T in G . Since we allow revisiting vertices, it is possible that $v_j = v_{j'}$ for some $j \neq j'$.

Now define \tilde{q} component-wise where \tilde{q}_j corresponds to the number of times the vertex j is visited in the walk W . Since the walk has a length T , trivially $\sum_j \tilde{q}_j = T$, satisfying (9e). Then, \tilde{y}, \tilde{q} can be defined

so that we have a feasible solution to the MP. Hence, if we now show that the CG defined by the nonzero components of \tilde{q} has a Hamiltonian path, then the corresponding \tilde{x} will be feasible for the AWT due to [Theorem 22](#).

Now, in the CG, construct the path $P = (v_1, n(v_1) + 1), (v_2, n(v_2) + 1), \dots, (v_T, n(v_T) + 1)$ where $n(v_j)$ records the number of times v_j has appeared in the path P earlier so far. We note that each term in the path P is indeed a vertex of the CG (which has T vertices) and that they are all visited exactly once, implying that P is the required Hamiltonian path. \blacksquare \square

The general mechanism involving the CP component is provided in [Algorithm 1](#), which is the entire algorithm we test in [Section 6](#). The CP component checks whether the solution returned by the CMP can be made valid in some way, i.e., if there exist solutions using only the configurations in \mathcal{A} (i.e., the configurations that appear in the optimal CMP solution, see [Definition 25](#)) with integer values, such that they do not violate the transition constraints. By providing these feasible solutions, it provides an upper bound to the AWT. As soon as a set of configurations is given to CP, a cut is added to the CMP, which eliminates all solutions to CMP which only consist of the configurations provided to CP. Namely, if \mathcal{A} indexes the configurations given to CP, we add a cut $\sum_{j \in \mathcal{A}} q_j \leq T - 1$, indicating that at least one of the configurations must be outside the set indexed by \mathcal{A} .

Let k' be the cardinality of \mathcal{A} . A CCG is associated with solution q — this CCG forms the basis for a deterministic finite automaton. This automaton A is defined by a tuple $(Q, \Sigma, \delta, q_0, F)$ of states Q , alphabet Σ , transition function $\delta : Q \times \Sigma \rightarrow Q$, initial state $q_0 \in Q$, and final states $F \subseteq Q$, with $Q = \{0, \dots, k'\}$, $\Sigma = \{1, \dots, k'\}$, $\delta = \{(u, v) \rightarrow v : (u, v) \in E\} \cup \{(0, u) \rightarrow u : u \in F\}$, $q_0 = 0$, and $F = \{1, \dots, k'\}$. In other words, this automaton accepts any valid sequence of k' configurations, with dummy configuration $q_0 = 0$ being the initial state. Let LB be the lower bound given by the CMP, and UB be an upper bound. Finally, let Ω be the collection of all index sets \mathcal{A}' for which the CP model has been previously solved. There are T decision variables z with domains $\{0, \dots, k' - 1\}$, with z_t indicating which CMP configuration from \mathcal{A} is used at time point t . These variables are constrained by

$$\min \phi \tag{12a}$$

$$\phi = \max(c) - \min(c) \tag{12b}$$

$$LB \leq \phi \leq UB - 1 \tag{12c}$$

$$\text{cost_regular}(z, A, \tau_{i*}, c_i) \quad i = 1, \dots, n \tag{12d}$$

$$\text{at_most}(T - 1, z, \omega) \quad \forall \omega \subset \mathcal{A} : \exists \mathcal{A}' \in \Omega, \omega \subseteq \mathcal{A}' \tag{12e}$$

$$z_t \in \mathcal{A} \quad t = 1, \dots, T \tag{12f}$$

$$c_i \in \mathbb{Z}_{\geq 0} \quad i = 1, \dots, n \tag{12g}$$

The objective [\(12a\)](#) is to minimize the unfairness, i.e., the difference between the zone which is covered the most, and that which is covered the least [\(12b\)](#). Any feasible solution should be strictly better than the upper bound, and search by the CP solver can be interrupted as soon as a feasible solution is found whose objective value coincides with LB [\(12c\)](#). By [Theorem 27](#), any sequence of configurations indexed by \mathcal{A} and of length T must correspond to a feasible solution of the AWT. This requirement is enforced by the **cost-regular** ([Demasse et al., 2006](#)) constraints [\(12d\)](#): A **cost-regular** constraint holds for zone i if z forms a word recognized by automaton A and if variable c_i is equal to the coverage of zone i over T : $c_i = \sum_{t=1}^T \tau_{i, z_t}$. Note that the CP model does not require that all configurations indexed by \mathcal{A} be used at least once: Since it checks all possible subsets of \mathcal{A} , the cut added to CMP does not remove any unexplored part of the search space. Finally, all subsets $\omega \subset \mathcal{A}$ previously explored, i.e. being as well a subset of some index set \mathcal{A}' in Ω , are considered by the model using the **at_most** constraint [\(12e\)](#): At most $T - 1$ variables in z can take on values in ω . The motivation for this is illustrated in [Example 28](#).

Example 28. Assume that $\mathcal{A} = \{j', j'', j'''\}$. If (12) had previously solved $\mathcal{A}' = \{j', j''\}$, this solution space would be explored again since the configurations in \mathcal{A} are not constrained to be used at least once. Adding constraint $\text{at_most}(T - 1, z, \{j', j''\})$ in the current iteration avoids this.

The CP solver we use is OR-Tools, which currently does not implement the `cost-regular` constraint. We thus replaced (12b) and (12d) with the equivalent but less efficient⁴ reformulation

$$\phi = \max_{i=1}^n \sum_{t=1}^T \tau_{i,z_t} - \min_{i=1}^n \sum_{t=1}^T \tau_{i,z_t} \\ \text{regular}(z, A).$$

5.4 The Final Algorithm

The general working of the final algorithm is presented formally in Algorithm 1.

In each iteration of the final algorithm, we solve a linear program and make a call to the constraint programming solver. The linear program is basically the continuous relaxation of the MP with any subsequent cutting planes (Step 6 in Algorithm 1). This is solved using column generation. The configurations which receive non-zero weights (indexed by \mathcal{A}) in the optimal LP solution are passed to the constraint programming solver to detect whether there exists a solution to the AWT that uses only configurations indexed by \mathcal{A} . Meanwhile, a cut is added to the linear program, that eliminates all solutions which use only the configurations indexed by \mathcal{A} .

With this, the large set of possible configurations are managed by the linear programming solver, which is powerful and efficient for large problems, while the CP-based solver only solves instances with a handful of configurations at a time.

6 Computational Experiments

In this section, we introduce the setting that we have used to evaluate computationally Algorithm 1 and we discuss the results of such an evaluation in the context of ambulance location and relocation. In particular, we use real data from the city of Utrecht, Netherlands, as well as synthetically-generated instances of varying sizes. Moreover, we consider a predetermined time horizon, i.e., a fixed T of size 30, with the understanding that a decision maker could reasonably make plans on a monthly basis.

Utrecht instance. We use the model defined in Section 5 to determine ambulance allocation in the city of Utrecht, Netherlands, using a dataset provided by the RIVM.⁵ The Utrecht instance contains 217 zones, 18 of which are bases. Since calls should be reached within 15 minutes, we consider that a zone is connected to another if an ambulance can travel between the two zones in under 15 minutes. This makes the graph about 28.4% dense. Moreover, the fleet of ambulance consists of 19 ambulance vehicles.

The efficiency measure we impose is that at least 95% of all the zones should be covered at all times. We consider that a zone is covered if sufficiently many ambulances are in the vicinity of that zone. In turn, we define the sufficient number of ambulances for a zone to be 1, 2, 3, or 4, based on the population density of the zone.

Synthetic instances. Reproducing the ratio of bases and edges to zones, we also generated synthetic instances⁶ of sizes 50, 100, 200, and 400 zones, using the Matérn cluster point process (Matérn, 1960), which helps in generating a distribution of zones mimicking a realistic urban setting. The number of ambulances for the instances is chosen to be just enough to ensure feasibility. The results are averaged over five instances of each size.

⁴In practice, the CP component of Algorithm 1 remains very fast, so this loss in efficiency is negligible.

⁵National Institute for Public Health and the Environment of the Netherlands.

⁶These instances can be found at <https://github.com/PhilippeOlivier/ambulances>.

Algorithm 1 The Final Algorithm

Input: The ambulance allocation problem with number of zones n , the bases \mathcal{B} , ζ_i , for $i = 1, \dots, n$, f and a_{ih} for $i, h = 1, \dots, n$, $r \in \mathbb{Z}_{\geq 0}$ and $T \in \mathbb{Z}_+$.

Output: $x(1), \dots, x(T)$ that is optimal to the AWT.

```
1:  $LB \leftarrow -\infty, UB \leftarrow +\infty, \mathcal{C} \leftarrow \emptyset, x^*(1), \dots, x^*(T) = NULL$ .
2: while  $UB > LB$  do
3:   Solve CMP, i.e., the continuous relaxation of the MP after adding each constraint in  $\mathcal{C}$ . Let  $q^*$  be
   the optimal solution and  $o^*$  be the optimal objective value.
4:    $LB \leftarrow o^*$ .
5:    $\mathcal{A} \leftarrow \{j : q_j^* > 0\}$ .
6:    $\mathcal{C} \leftarrow \mathcal{C} \cup \left\{ \sum_{j \in \mathcal{A}} q_j \leq T - 1 \right\}$ .
7:    $(x^\dagger(1), \dots, x^\dagger(T)), o^\dagger \leftarrow \text{CONSTRAINTPROGRAMMING}(q^*, n, a, T, LB, UB, \mathcal{C})$ .
8:   if  $o^\dagger < UB$  then
9:      $UB \leftarrow o^\dagger$  and  $(x^*(1), \dots, x^*(T)) \leftarrow (x^\dagger(1), \dots, x^\dagger(T))$ .
10:  end if
11: end while
12: return  $x^*(1), \dots, x^*(T)$ 
13:
14: function  $\text{CONSTRAINTPROGRAMMING}(q, n, a, T, LB, UB, \mathcal{C})$ 
15:    $G \leftarrow \text{CCG defined by } \mathcal{A}$ .
16:   Solve (12) on the graph  $G$  with appropriate values of  $\tau, c$ .
17:   if (12) is infeasible then
18:     return  $NULL, +\infty$ 
19:   else
20:     return  $(x^\dagger(1), \dots, x^\dagger(T)), o^\dagger$ .
21:   end if
22: end function
```

Testing environment and software. Testing was performed on an Intel 2.8 GHz processor with 8 GBs of RAM running Arch Linux, and the models were solved with Gurobi 9.0.1 and OR-Tools 8.0. A time limit of 1,200 seconds was imposed to solve each instance.

Results. It is easy to see that there are two parameters that are likely to influence the difficulty of the solving process: the size of the instance (number of zones) and the transition constraints, i.e., the flexibility we allow for relocating ambulances between zones on a daily basis. We would typically expect the size of the instance to be a significant factor, but this is likely to be mitigated by the column generation approach, which tends to scale well. The effects that the transition constraints will have on the solving process, however, are less clear. In order to assess such an effect, for each instance, we vary the number of ambulances that can shift bases on consecutive days (r in constraints (8k)) from $0.1m$ to m in increments of $0.1m$, where m is the total number of ambulances in the instance. We call this ratio maximum transition, MT . For each of these cases, we record the average of the time taken to solve these instances to optimality (capped at 1,200 CPU seconds). We record the final relative gap left for the instance, which is defined by $\frac{|UB_f - LB_f|}{UB_f}$, where LB_f and UB_f represent the values of the best lower and upper bounds at the time limit, respectively. We also record the initial relative gap for the instance, which is defined by $\frac{|UB_i - LB_i|}{UB_i}$, where LB_i represents the value of the initial lower bound (without any cuts), and UB_i represents the value of the initial, trivial upper bound. We present the final relative gaps associated with different classes of instances in Figure 2.

One can immediately observe from Figure 2 that the transition constraints are affecting the difficulty of the instances for $MT \leq 0.5$, while everything can be solved to optimality within the time limit for $MT \geq 0.5$

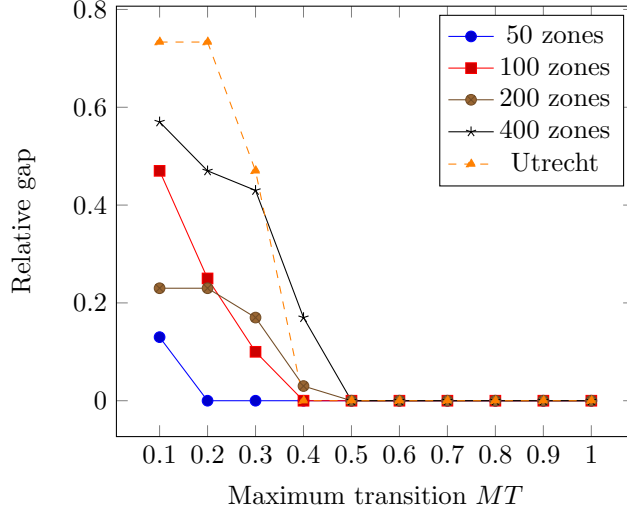


Figure 2: Final relative gaps with respect to the maximum transition MT . Gaps in the synthetic instances are averaged over 5 instances. Time limit of 1,200 CPU seconds.

independently of the number of zones. The larger MT the smaller the final gap, i.e., when there is more freedom in moving the ambulances around on a daily basis, [Algorithm 1](#) becomes very effective in computing optimal solutions.

We observe that for varying values of MT , the Utrecht instance has a higher relative gap than synthetic instances of comparable and even larger size. This discrepancy is the result of the distribution of the population densities across the zones. In the synthetic instances, the population densities are randomly distributed among the zones. The Utrecht instance, in contrast, exhibits several distinct clusters of varying sizes and population density distributions.⁷

Unsurprisingly, [Figure 2](#) also suggests that larger instances are more difficult because, typically, the relative gap is large when time limits are hit. This is confirmed in more details from the results in [Table 3](#), where we report, for every group of instances and every MT value, the initial gap (“i.gap”), the final gap at the time limit (“f.gap”), the number of generated columns (“#cols”), the number of solved instances (“#solved”),⁸ and the average time for the instances solved to optimality (“time”).

The results in [Table 3](#) show that for $MT > 0.7$ constraints (8k) are not binding, i.e., they do not cut off the obtained optimal solution to the rest of the model, and [Algorithm 1](#) behaves exactly like for $MT = 0.7$. The average time for the instances solved to optimality is often low, thus indicating that if we manage to prove optimality by completely closing the gap, then this is achieved quickly, generally with a few calls to `CONSTRAINTPROGRAMMING`. In contrast, if we do not succeed in closing the gap early, then it is unlikely to be closed at all in the end. In the cases where optimality is not proven, improvements are still made during the solving process mostly because of improvements on the upper bound value as `CONSTRAINTPROGRAMMING` tests more and more combinations of configurations. This suggests a few things. First, with a high enough MT , the initial lower bound is generally optimal, and feasible solutions whose objective values coincide with this bound can readily be found. Second, when the MT value is low, solutions close to the initial LB are harder to find, and the gap is harder to close. Third, the computational power of [Algorithm 1](#) is associated with the effectiveness in searching for primal solutions (upper bound improvement) versus the progress on the dual side (lower bound improvement), which appears to be very difficult.

[Tables 4 to 6](#) show the actual values of the solutions, as well as the associated coverages (the sum of all the

⁷Constructing these sophisticated clusters is not a trivial task, which is why we settled on a random distribution of population densities for the synthetic instances. By randomizing the placement of the population for the Utrecht instance, its gap becomes similar to that of the synthetic instances.

⁸The entry “#solved” takes an integer value between 0 and 5 for synthetic instances and 0/1 for the Utrecht instance.

MT	50 zones					100 zones					200 zones				
	i.gap	f.gap	#cols	#solved	time	i.gap	f.gap	#cols	#solved	time	i.gap	f.gap	#cols	#solved	time
0.1	0.13	0.13	738	4	0.3	0.55	0.47	1120	1	0.2	0.23	0.23	714	3	2.5
0.2	0.13	0.00	3	5	0.4	0.55	0.25	674	2	1.2	0.23	0.23	710	3	2.5
0.3	0.13	0.00	3	5	0.4	0.55	0.10	11	4	2.5	0.23	0.17	537	3	2.6
0.4	0.13	0.00	3	5	0.4	0.55	0.00	9	5	2.7	0.23	0.02	108	4	82.7
0.5	0.13	0.00	3	5	0.4	0.55	0.00	9	5	2.3	0.23	0.00	19	5	7.0
0.6	0.13	0.00	3	5	0.4	0.55	0.00	9	5	2.1	0.23	0.00	15	5	7.2
0.7	0.13	0.00	3	5	0.4	0.55	0.00	9	5	2.1	0.23	0.00	15	5	6.7
0.8	0.13	0.00	3	5	0.4	0.55	0.00	9	5	2.2	0.23	0.00	15	5	6.0
0.9	0.13	0.00	3	5	0.4	0.55	0.00	9	5	2.2	0.23	0.00	15	5	6.0
1	0.13	0.00	3	5	0.4	0.55	0.00	9	5	2.1	0.23	0.00	15	5	6.0

MT	400 zones					Utrecht				
	i.gap	f.gap	#cols	#solved	time	i.gap	f.gap	#cols	#solved	time
0.1	0.57	0.57	248	1	2.7	0.73	0.73	786	0	-
0.2	0.57	0.47	246	1	2.7	0.73	0.73	803	0	-
0.3	0.57	0.43	238	1	2.7	0.73	0.47	757	0	-
0.4	0.57	0.17	174	3	147.4	0.73	0.00	353	1	315.6
0.5	0.57	0.00	93	5	273.6	0.73	0.00	252	1	122.3
0.6	0.57	0.00	65	5	63.6	0.73	0.00	96	1	49.6
0.7	0.57	0.00	60	5	49.6	0.73	0.00	96	1	50.3
0.8	0.57	0.00	60	5	47.1	0.73	0.00	96	1	48.4
0.9	0.57	0.00	60	5	47.4	0.73	0.00	96	1	49.8
1	0.57	0.00	60	5	49.1	0.73	0.00	96	1	54.7

Table 3: Detailed results for Algorithm 1 on both synthetic and real-world Utrecht instances. Time limit of 1,200 CPU seconds. The results of the synthetic instances are averaged over 5 instances.

MT	50 zones					100 zones					200 zones					400 zones					Utrecht				
	g	h	$g-h$	cov		g	h	$g-h$	cov		g	h	$g-h$	cov		g	h	$g-h$	cov		g	h	$g-h$	cov	
0.1	30	24	6	99.20%	30	10	20	96.20%	30	6	24	97.33%	30	6	24	98.45%	-	-	-	-	-	-	-	-	-
0.2	30	28	2	99.47%	30	16	14	96.20%	30	6	24	97.33%	30	12	18	98.45%	-	-	-	-	-	-	-	-	-
0.3	30	28	2	99.47%	30	20	10	96.53%	30	9	21	97.33%	30	16	14	98.45%	30	15	15	97.23%	30	15	15	97.23%	30
0.4	30	28	2	99.47%	30	23	7	96.43%	30	13	17	97.33%	30	22	8	98.19%	30	22	8	95.39%	30	22	8	95.39%	30
0.5	30	28	2	99.47%	30	23	7	96.43%	30	13	17	97.15%	30	23	7	97.23%	30	22	8	95.48%	30	22	8	95.48%	30
0.6	30	28	2	99.47%	30	23	7	96.43%	30	13	17	97.35%	30	23	7	97.21%	30	22	8	95.39%	30	22	8	95.39%	30
0.7	30	28	2	99.47%	30	23	7	96.43%	30	13	17	97.00%	30	23	7	97.23%	30	22	8	95.39%	30	22	8	95.39%	30
0.8	30	28	2	99.47%	30	23	7	96.43%	30	13	17	97.22%	30	23	7	97.13%	30	22	8	95.39%	30	22	8	95.39%	30
0.9	30	28	2	99.47%	30	23	7	96.43%	30	13	17	97.22%	30	23	7	97.23%	30	22	8	95.39%	30	22	8	95.39%	30
1	30	28	2	99.47%	30	23	7	96.43%	30	13	17	97.22%	30	23	7	97.17%	30	22	8	95.39%	30	22	8	95.39%	30

Table 4: Enforcing coverage constraints at 95%. Solution values for both synthetic and real-world Utrecht instances: g , h , $g-h$ (defined in (9)), and “cov” (the average coverage, i.e., the efficiency). The results of the synthetic instances are averaged over 5 instances.

MT	50 zones					100 zones					200 zones					400 zones					Utrecht				
	g	h	$g-h$	cov		g	h	$g-h$	cov		g	h	$g-h$	cov		g	h	$g-h$	cov		g	h	$g-h$	cov	
0.1	30	24	6	98.80%	30	15	15	95.27%	30	6	24	94.40%	30	9	21	98.45%	-	-	-	-	-	-	-	-	-
0.2	30	29	1	98.67%	30	23	7	95.10%	30	6	24	94.40%	30	13	17	98.45%	-	-	-	-	-	-	-	-	-
0.3	30	29	1	98.67%	30	23	7	94.17%	30	12	18	94.40%	30	15	15	98.45%	30	19	11	90.32%	30	19	11	90.32%	30
0.4	30	29	1	98.67%	30	24	6	94.17%	30	20	10	94.30%	30	23	7	97.70%	30	22	8	90.32%	30	22	8	90.32%	30
0.5	30	29	1	98.67%	30	24	6	94.17%	30	21	9	93.33%	30	25	5	95.31%	30	23	7	91.90%	30	23	7	91.90%	30
0.6	30	29	1	98.67%	30	24	6	94.17%	30	21	9	93.65%	30	25	5	95.42%	30	23	7	91.84%	30	23	7	91.84%	30
0.7	30	29	1	98.67%	30	24	6	94.17%	30	21	9	93.30%	30	25	5	95.19%	30	23	7	91.53%	30	23	7	91.53%	30
0.8	30	29	1	98.67%	30	24	6	94.17%	30	21	9	93.30%	30	25	5	95.05%	30	23	7	91.97%	30	23	7	91.97%	30
0.9	30	29	1	98.67%	30	24	6	94.17%	30	21	9	93.28%	30	25	5	95.10%	30	23	7	92.00%	30	23	7	92.00%	30
1	30	29	1	98.67%	30	24	6	94.17%	30	21	9	93.30%	30	25	5	95.07%	30	23	7	92.00%	30	23	7	92.00%	30

Table 5: Enforcing coverage constraints at 90%. Solution values for both synthetic and real-world Utrecht instances: g , h , $g-h$ (defined in (9)), and “cov” (the average coverage, i.e., the efficiency). The results of the synthetic instances are averaged over 5 instances.

MT	50 zones					100 zones					200 zones					400 zones					Utrecht				
	g	h	$g-h$	cov		g	h	$g-h$	cov		g	h	$g-h$	cov		g	h	$g-h$	cov		g	h	$g-h$	cov	
0.1	30	24	6	97.60%	30	10	20	93.60%	30	6	24	5616%	30	10	20	96.45%	-	-	-	-	-	-	-	-	-
0.2	30	29	1	98.40%	30	24	6	92.57%	30	6	24	5616%	30	13	17	96.45%	-	-	-	-	-	-	-	-	-
0.3	30	29	1	98.40%	30	25	5	92.90%	30	14	16	5616%	30	17	13	96.45%	30	22	8	94.47%	30	22	8	94.47%	30
0.4	30	29	1	98.40%	30	25	5	92.57%	30	20	10	5610%	30	22	8	96.45%	30	25	5	90.48%	30	25	5	90.48%	30
0.5	30	29	1	98.40%	30	25	5	92.57%	30	22	8	5477%	30	24	6	92.53%	30	25	5	90.48%	30	25	5	90.48%	30
0.6	30	29	1	98.40%	30	25	5	92.57%	30	22	8	5464%	30	25	5	92.52%	30	25	5	90.48%	30	25	5	90.48%	30
0.7	30	29	1	98.40%	30	25	5	92.57%	30	22	8	5448%	30	25	5	93.19%	30	25	5	90.48%	30	25	5	90.48%	30
0.8	30	29	1	98.40%	30	25	5	92.57%	30	22	8	5454%	30	25	5	93.14%	30	25	5	90.48%	30	25	5	90.48%	30
0.9	30	29	1	98.40%	30	25	5	92.57%	30	22	8	5453%	30	25	5	93.19%	30	25	5	90.48%	30	25	5	90.48%	30
1	30	29	1	98.40%	30	25	5	92.57%	30	22	8	5453%	30	25	5	93.14%	30	25	5	90.48%	30	25	5	90.48%	30

Table 6: Enforcing coverage constraints at 85%. Solution values for both synthetic and real-world Utrecht instances: g , h , $g-h$ (defined in (9)), and “cov” (the average coverage, i.e., the efficiency). The results of the synthetic instances are averaged over 5 instances.

zones’ coverages over the time horizon, i.e., the efficiency). In Table 4 the coverage constraints are enforced

at 95%, i.e., all the configurations must cover at least 95% of the zones. In Table 5 those constraints are enforced at 90%, and in Table 6 at 85%. Comparing the three tables, we can see a clear correlation between fairness and efficiency in virtually all instances: Enforcing a higher coverage decreases the fairness of the solutions.

Finally, we are also interested in identifying the time spent in the column generation part as opposed to the CONSTRAINTPROGRAMMING routine. To this end, we track the number of calls made to the function CONSTRAINTPROGRAMMING and also measure the total time spent in calls to the function. The number of generated columns and the number of calls to CONSTRAINTPROGRAMMING are detailed in Table 7.

	50 zones		100 zones		200 zones		400 zones		Utrecht	
<i>MT</i>	#cols	#calls	#cols	#calls	#cols	#calls	#cols	#calls	#cols	#calls
0.1	738	736	1120	1108	714	608	248	72	786	210
0.2	3	1	674	662	710	605	246	72	803	217
0.3	3	1	11	3	537	435	238	66	757	202
0.4	3	1	9	1	108	58	174	22	353	54
0.5	3	1	9	1	19	3	93	6	252	21
0.6	3	1	9	1	15	1	65	2	96	1
0.7	3	1	9	1	15	1	60	1	96	1

Table 7: Number of columns and calls to CONSTRAINTPROGRAMMING associated with the maximum transition (MT) constraints of the various instance sizes.

The results in Table 7 show that when enough freedom in the day-to-day movement of ambulances is allowed, optimality can generally be proven immediately, with just one call to the CONSTRAINTPROGRAMMING routine. We also notice that the ratio of columns to the number of calls to the routine increases with the instance size. Such a ratio for the Utrecht instance is disproportionately high, due again to the distribution of the population densities.

Finally, Table 8 shows the average amount of time spent generating columns (“CG”), that spent enforcing the transition constraints through CONSTRAINTPROGRAMMING (“CP”) and the ratio between these two CPU times (“ratio”) for the most difficult case, i.e., $MT = 0.1$.

Instance	time		ratio
	CG	CP	
50 zones	86	154	0.56
100 zones	411	550	0.75
200 zones	331	151	2.19
400 zones	711	253	2.81
Utrecht	1071	129	7.96

Table 8: CPU times in seconds spent in column generation and CONSTRAINTPROGRAMMING, with $MT = 0.1$.

7 Conclusion

We introduced an abstract framework for solving a sequence of fair allocation problems, such that fairness is achieved over time. For some relevant special cases, we have been able to give theoretical proofs for the time horizon required for perfect fairness. We described a general integer programming formulation for this problem, as well as a formulation based on column generation and constraint programming. This latter formulation can be used in a practical context, as shown by the ambulance location problem applied to the city of Utrecht. The largest synthetic instances suggest that this formulation would scale well to regions twice the size of Utrecht, given that the freedom of movement of the ambulances is not overly restricted.

Acknowledgements

The authors would like to thank the National Institute for Public Health and the Environment of the Netherlands (RIVM) for access to the data of their ambulance service. We are grateful to the anonymous referees for useful comments and remarks.

References

- Martin Aleksandrov and Toby Walsh. Online fair division: A survey. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 13557–13562, 2020.
- Ahmet Alkan, Gabrielle Demange, and David Gale. Fair allocation of indivisible goods and criteria of justice. *Econometrica: Journal of the Econometric Society*, pages 1023–1039, 1991.
- Egon Balas and Robert Jeroslow. Canonical cuts on the unit hypercube. *SIAM Journal on Applied Mathematics*, 23(1):61–69, 1972.
- Evrpidis Bampis, Bruno Escoffier, and Sasa Mladenovic. Fair resource allocation over time. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, AAMAS ’18, page 766–773, Richland, SC, 2018. International Foundation for Autonomous Agents and Multiagent Systems.
- Luce Brotcorne, Gilbert Laporte, and Frédéric Semet. Ambulance location and relocation models. *European Journal of Operational Research*, 147(3):451–463, 2003. ISSN 0377-2217. doi: [https://doi.org/10.1016/S0377-2217\(02\)00364-8](https://doi.org/10.1016/S0377-2217(02)00364-8). URL <http://www.sciencedirect.com/science/article/pii/S0377221702003648>.
- Oscar R. Burt and Curtis C. Harris. Letter to the editor—apportionment of the u.s. house of representatives: A minimum range, integer solution, allocation problem. *Operations Research*, 11(4):648–652, 1963. doi: 10.1287/opre.11.4.648. URL <https://doi.org/10.1287/opre.11.4.648>.
- Claudia D’Ambrosio, Antonio Frangioni, Leo Liberti, and Andrea Lodi. On interval-subgradient and no-good cuts. *Operations Research Letters*, 38(5):341–345, 2010.
- Sophie Demasse, Gilles Pesant, and Louis-Martin Rousseau. A cost-regular based hybrid column generation approach. *Constraints*, 11(4):315–333, Dec 2006. ISSN 1572-9354. doi: 10.1007/s10601-006-9003-7. URL <https://doi.org/10.1007/s10601-006-9003-7>.
- Ohad Eisenhandler and Michal Tzur. The humanitarian pickup and distribution problem. *Operations Research*, 67(1):10–32, 2019. doi: 10.1287/opre.2018.1751. URL <https://doi.org/10.1287/opre.2018.1751>.
- Gerald Gamrath, Daniel Anderson, Ksenia Bestuzheva, Wei-Kun Chen, Leon Eifler, Maxime Gasse, Patrick Gemander, Ambros Gleixner, Leona Gottwald, Katrin Halbig, Gregor Hendel, Christopher Hojny, Thorsten Koch, Pierre Le Bodic, Stephen J. Maher, Frederic Matter, Matthias Miltenberger, Erik Mühmer, Benjamin Müller, Marc E. Pfetsch, Franziska Schlösser, Felipe Serrano, Yuji Shinano, Christine Tawfik, Stefan Vigerske, Fabian Wegscheider, Dieter Weninger, and Jakob Witzig. The SCIP Optimization Suite 7.0. ZIB-Report 20-10, Zuse Institute Berlin, March 2020a. URL <http://nbn-resolving.de/urn:nbn:de:0297-zib-78023>.
- Gerald Gamrath, Daniel Anderson, Ksenia Bestuzheva, Wei-Kun Chen, Leon Eifler, Maxime Gasse, Patrick Gemander, Ambros Gleixner, Leona Gottwald, Katrin Halbig, Gregor Hendel, Christopher Hojny, Thorsten Koch, Pierre Le Bodic, Stephen J. Maher, Frederic Matter, Matthias Miltenberger, Erik Mühmer, Benjamin Müller, Marc E. Pfetsch, Franziska Schlösser, Felipe Serrano, Yuji Shinano, Christine Tawfik, Stefan Vigerske, Fabian Wegscheider, Dieter Weninger, and Jakob Witzig. The SCIP Optimization Suite 7.0. Technical report, Optimization Online, March 2020b. URL http://www.optimization-online.org/DB_HTML/2020/03/7705.html.

- Michel Gendreau, Gilbert Laporte, and Frédéric Semet. A dynamic model and parallel tabu search heuristic for real-time ambulance relocation. *Parallel Computing*, 27(12):1641–1653, 2001. ISSN 0167-8191. doi: [https://doi.org/10.1016/S0167-8191\(01\)00103-X](https://doi.org/10.1016/S0167-8191(01)00103-X). URL <http://www.sciencedirect.com/science/article/pii/S016781910100103X>. Applications of parallel computing in transportation.
- Jessica L. Heier Stamm, Nicoleta Serban, Julie Swann, and Pascale Wortley. Quantifying and explaining accessibility with application to the 2009 h1n1 vaccination campaign. *Health Care Management Science*, 20(1):76–93, Mar 2017. ISSN 1572-9389. doi: 10.1007/s10729-015-9338-y. URL <https://doi.org/10.1007/s10729-015-9338-y>.
- Stephen Jacobsen. On marginal allocation in single constraint min-max problems. *Management Science*, 17(11):780–783, 1971. doi: 10.1287/mnsc.17.11.780. URL <https://doi.org/10.1287/mnsc.17.11.780>.
- Naoki Katoh, Toshihide Ibaraki, and H. Mine. An algorithm for the equipollent resource allocation problem. *Mathematics of Operations Research*, 10(1):44–53, 1985. URL <https://EconPapers.repec.org/RePEc:inm:ormoor:v:10:y:1985:i:1:p:44-53>.
- Bernard O. Koopman. The optimum distribution of effort. *Journal of the Operations Research Society of America*, 1(2):52–63, 1953. doi: 10.1287/opre.1.2.52. URL <https://doi.org/10.1287/opre.1.2.52>.
- Hanan Luss. *Equitable Resource Allocation: Models, Algorithms and Applications*. Information and Communication Technology Series,. Wiley, 2012. ISBN 9781118449219. URL https://books.google.ca/books?id=Z2_3oWjnASkC.
- François Margot. Symmetry in integer linear programming. In *50 Years of Integer Programming 1958-2008*, pages 647–686. Springer, 2010.
- Bertil Matérn. Spatial variation – stochastic models and their applications to some problems in forest survey sampling investigations. *Report of the Forest Research Institute of Sweden*, 49(5):1–144, 1960.
- Włodzimierz Ogryczak, Hanan Luss, Michał Pióro, Dritan Nace, and Artur Tomaszewski. Fair optimization and networks: A survey. *Journal of Applied Mathematics*, 2014:612018, Sep 2014. ISSN 1110-757X. doi: 10.1155/2014/612018. URL <https://doi.org/10.1155/2014/612018>.
- Ogryczak, Włodzimierz. Fair optimization – methodological foundations of fairness in network resource allocation. In *2014 IEEE 38th International Computer Software and Applications Conference Workshops*, pages 43–48, 2014.
- Evan L. Porteus and Jonathan S. Yormark. More on min-max allocation. *Management Science*, 18(9): 502–507, 1972. doi: 10.1287/mnsc.18.9.502. URL <https://doi.org/10.1287/mnsc.18.9.502>.
- Ariel D Procaccia. Cake cutting algorithms. In *Handbook of computational social choice, chapter 13*. Citeseer, 2015.
- Naoki Katoh Toshihide Ibaraki. *Resource Allocation Problems: Algorithmic Approaches*, page 1. Foundations of Computing. The MIT Press, 1988. ISBN 0262090279,9780262090278.
- Zeev Zeitlin. Minimization of maximum absolute deviation in integers. *Discrete Applied Mathematics*, 3(3):203 – 220, 1981. ISSN 0166-218X. doi: [https://doi.org/10.1016/0166-218X\(81\)90017-2](https://doi.org/10.1016/0166-218X(81)90017-2). URL <http://www.sciencedirect.com/science/article/pii/0166218X81900172>.

Appendix

A Dual of the CMP

$$\max \quad T\mu \quad (13a)$$

$$\text{s.t.} \quad \sum_{i=1}^n \alpha_i = 1 \quad (g) \quad (13b)$$

$$\sum_{i=1}^n \beta_i = 1 \quad (h) \quad (13c)$$

$$\lambda_i - \alpha_i + \beta_i = 0 \quad (y_i) \quad \forall i = 1, \dots, n \quad (13d)$$

$$\mu \leq \frac{1}{T} \sum_{i=1}^n \tau_{ij} \lambda_i \quad (q_j) \quad \forall j = 1, \dots, k \quad (13e)$$

$$\alpha_i, \beta_i \geq 0 \quad \forall i = 1, \dots, n. \quad (13f)$$