

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

Andrea Lodi^{1,2}, Philippe Olivier^{1,2}, Gilles Pesant², and Sriram Sankaranarayanan^{*3}

¹*CERC*

²*Polytechnique Montréal, Canada*

³*IIM 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

*srirams@iima.ac.in

39 efficiency and fairness in resource allocation is common in communication networks (Ogryczak et al., 2014,
40 Ogryczak, Włodzimierz, 2014), and can be found in various other fields, such as in some pickup and delivery
41 problems (Eisenhandler and Tzur, 2019).

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

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

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

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

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

73 2 Preliminaries

74 In this section, we formally define the terms used throughout the manuscript. We study resource allocation
75 in the context of fairness over time, meaning that resources are allocated to stakeholders over some time
76 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
77 stakeholders among whom the allocations should be done in a fair way.

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

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

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

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

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

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

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

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

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

90 Now, in this section, we first define the properties that a general measure of unfairness should have for it
91 to be considered valid in our context. To start with, we state that (i) if all the stakeholders of the problem
92 get the same benefit or utility, then the unfairness associated with such an allocation of benefits should be
93 0, and (ii) the unfairness associated with benefits should be invariant to permutations of the ordering of
94 stakeholders.

95 **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
96 the allocations, if the i -th stakeholder gets a benefit of y_i . Such a function ϕ satisfies the following:*

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

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

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

102 **Definition 3** (Inefficiency function). *Given \mathcal{X} and the benefit function τ , the inefficiency function $\eta : \mathcal{X} \rightarrow$
103 $[0, 1]$ is defined as*

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

104 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.

105 *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
106 (or allocation) that maximizes the sum of benefits, $\eta(x) = 0$, while for the least efficient x , we have $\eta(x) = 1$.
107 Thus, η serves as a method of normalization of the objective values.

108 **Definition 5.** *Equivalence classes defined due to the function η and partitioning \mathcal{X} are, for any $\bar{\eta} \in [0, 1]$*

$$\tilde{\mathcal{X}}_{\bar{\eta}} := \{x \in \mathcal{X} : \eta(x) = \bar{\eta}\}. \tag{4}$$

109 The above definition enables to succinctly represent a target to find fair solutions *only* if they are also
110 sufficiently efficient. For instance, one might restrict the feasible set to $\cup_{\eta: \eta < \bar{\eta}} \tilde{\mathcal{X}}_{\eta}$.

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

Definition 6 (Single-period fair allocation (SPFA) problem). *Given $\bar{\eta} \in [0, 1]$, the single-period fair allocation problem is to solve*

$$\min_x \phi(\tau(x)) \quad (5a)$$

$$s.t. \quad x \in \bigcup_{\eta: \eta \in [0, \bar{\eta}]} \tilde{\mathcal{X}}_\eta. \quad (5b)$$

113 **Example 7** motivates and provides a mean of validation for model (5). It shows that with reasonable
114 choices, we retrieve the intuitively fair solution.

115 **Example 7.** *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
116 that $[\tau(x)]_i = \tau_i x_i$ for some scalars $\tau_i \geq 0$, 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

$$x_i^* = \frac{g}{\tau_i}, \quad \forall i \quad (6a)$$

$$\text{where } g = \frac{1}{\sum_{i=1}^n \frac{1}{\tau_i}}. \quad (6b)$$

117 *Clearly, each stakeholder enjoys a benefit of g , and hence the unfairness associated with this allocation is 0.*

118 *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

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

119 We now define the fairness-over-time problem.

Definition 8 (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*

$$\min_{x(t), y} \phi(y) \quad (7a)$$

$$s.t. \quad y = \frac{1}{T} \sum_{t=1}^T \tau(x(t)) \quad (7b)$$

$$x(t) \in \bigcup_{\eta: \eta \in [0, \bar{\eta}]} \tilde{\mathcal{X}}_\eta, \quad \forall i = 1, \dots, T. \quad (7c)$$

120

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

123 **Example 9** (Usefulness of the T -PFA). *Let $\mathcal{X} = \{x \in \{0, 1\}^2 : x_1 + x_2 = 1\}$. Further, let $\tau(x) = (2x_1, x_2)$.
124 Let $\bar{\eta} = 1$. Note that, in the case of the SPFA, the optimal objective is necessarily nonzero, since the
125 benefits of the two stakeholders are unequal for every feasible solution. However, consider the 3-PFA. Now,
126 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
127 is 0, which is strictly better than the SPFA solution.*

3 Simple Cases

3.1 Convex \mathcal{X}

We have motivated in [Example 9](#) that ensuring fairness over multiple rounds could offer better results than seeking fairness on a per-round basis. We now show that this holds only for a nonconvex \mathcal{X} . In other words, if \mathcal{X} is convex, we necessarily get no improvement in T -PFA over SPFA.

Theorem 10. *Let \mathcal{X} be convex. Further, let τ be a linear function. Let $\bar{\eta}$ be given. Let f^* be the optimal value of the SPFA problem and let f_2^*, f_3^*, \dots be the optimal values of the T -PFA problem with $T = 2, 3, \dots$. Then $f^* = f_2^* = f_3^* \dots$*

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 $\hat{\mathcal{X}} = \cup_{\eta: \eta \in [0, \bar{\eta}]} \mathcal{X}_\eta$. 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, we show that it is not necessarily due to the fact that perfect fairness is obtainable in SPFA. [Example 11](#) 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 10](#), perfect fairness is neither possible with multiple rounds of allocation.

Example 11. *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

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

for some $\rho \in \mathbb{R}_+$. Solving this linear system gives the unique fair allocation as allocations of the form $(x_1, x_2, x_3) = (-0.5\rho, \rho, \rho)$, 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}

The first part of [Theorem 12](#) states that for any n and a \mathcal{X} of the specified form, perfect fairness is possible precisely after \bar{T} periods, 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.

Theorem 12. *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 $L = \text{LCM}(\tau_1, \dots, \tau_n)$. Then,*

- 159 1. Perfect fairness cannot be obtained in T periods where T is strictly lesser than $\bar{T} = \left\lceil \frac{1}{a} \sum_{i=1}^n \frac{L}{\tau_i} \right\rceil$.
 160 However, if $1 > \bar{\eta} \geq \frac{\tau_1 - \min(\tau_n, \tau_1/a)}{\tau_1}$, and if $\bar{T} > 1$, perfect fairness can be achieved in exactly \bar{T}
 161 periods.
- 162 2. Given any $n \geq 2$, there exist $\tau_1 \geq \tau_2 \geq \dots \geq \tau_n$ and $a \in \mathbb{Z}_{>0}$ such that perfect fairness is unattainable
 163 for T -PFA where $T \leq \bar{T}$ as above with $\bar{\eta} < \frac{\tau_1 - \min(\tau_n, \tau_1/a)}{\tau_1}$.

Proof. For the first part, observe that an unfairness of 0 can be achieved if and only if the total benefit over time for each stakeholder should add up to the same value. Since x_i and τ_i are all integers and we have integer number of periods, the smallest value that the benefits can all add up to is 0, and the next smallest value is L . But the strict inequality $\bar{\eta} < 1$ excludes the trivial allocation of $(0, 0, \dots, 0)$. A total benefit of L can be achieved for stakeholder i if and only if $\sum_{t=1}^T x_i(t) = L/\tau_i$. Summing this over all stakeholders, we get

$$\begin{aligned}
 \sum_{i=1}^n \frac{L}{\tau_i} &= \sum_{i=1}^n \sum_{t=1}^T x_i(t) \\
 &= \sum_{t=1}^T \sum_{i=1}^n x_i(t) \\
 &\leq \sum_{t=1}^T a \\
 &= aT \\
 \implies \frac{1}{a} \sum_{i=1}^n \frac{L}{\tau_i} &\leq T
 \end{aligned}$$

164 and the ceiling follows from the fact that T is an integer. Now that we have shown that we need at least
 165 \bar{T} -PFA, we show that the unfairness can be made to 0 in exactly \bar{T} periods. For this, consider a reverse-
 166 lexicographic allocation method where all allocations are made to stakeholder n until the total allocation
 167 sums to L/τ_n . Next, allocate similarly for $n-1$, $n-2$ and so on up to stakeholder 1. Observe that each
 168 of these allocations at any time t is in \mathcal{X} by construction. To show that the allocation in each step has an
 169 inefficiency value at most $\bar{\eta}$, consider the two most inefficient allocations. This could happen in the first
 170 period, when all allocations are to the stakeholder n , i.e., $x(1) = (0, 0, \dots, 0, a)$. The inefficiency function
 171 value for 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, the most inefficient allocation could be in
 172 the last period, where a single allocation is made to the first stakeholder, i.e., $x(T) = (1, 0, 0, \dots, 0)$. The
 173 inefficiency function value for this allocation is $\frac{a\tau_1 - \tau_1}{a\tau_1 - 0} = \frac{\tau_1 - \tau_1/a}{\tau_1} \leq \bar{\eta}$. This proves the first part.

We prove the second part for $n=2$ first, and then generalize it to $n \geq 3$. Let the target $\bar{T} = \hat{T} \geq 2$. Choose $\tau_1 = (\hat{T}-1)a$ and $\tau_2 = 1$. The LCM $L = (\hat{T}-1)a$. From the first part, perfect fairness is not obtainable for $T < \frac{1}{a} \left(\frac{(\hat{T}-1)a}{(\hat{T}-1)a} + \frac{(\hat{T}-1)a}{1} \right) = (\hat{T}-1) + \frac{1}{a}$ periods. But from the integrality of the number of periods, perfect fairness is not possible using up to $\hat{T}-1$ periods. We now show that perfect fairness is not possible with \hat{T} periods either. Note that the set of solutions where perfect fairness is achieved in \hat{T} time steps are all in the form

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

$$x(\widehat{T}) = \left(\sum_{i=1}^{\widehat{T}-1} \delta_i, 1 \right)$$

174 and its permutations over time for some integers $\delta_i \in \{1, 2, \dots, a\}$ such that $0 \leq \sum_{i=1}^{\widehat{T}-1} \delta_i \leq a - 1$. Now
 175 consider the inefficiency function of the allocation $x(1)$; $\eta(x(1)) = \frac{(\widehat{T}-1)a-1+\delta_1/a}{\widehat{T}-1} > \bar{\eta}$ shows that it is
 176 infeasible. Thus, any perfectly fair allocation for \bar{T} -PFA problem is an infeasible allocation, providing the
 177 necessary counterexample. To extend this for $n > 2$, choose $\tau_2 = \tau_3 = \dots = \tau_{n-1} = 2$, and analogous
 178 arguments provide the counterexample. ■ □

179 **Example 13.** Consider [Theorem 12](#) applied to [Example 9](#). We have $\mathcal{X} = \{(1, 0), (0, 1)\}$, and thus $a = 1$.
 180 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, $L = LCM(2, 1) = 2$.
 181 Then, $\bar{T} = \left\lceil \frac{1}{a} \sum_{i=1}^n \frac{L}{\tau_i} \right\rceil = \left\lceil \frac{1}{1} \left(\frac{2}{2} + \frac{2}{1} \right) \right\rceil = 3$.

182 3.3 Sparse simplicial \mathcal{X}

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

187 **Theorem 14.** 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
 188 counts the number of non-zero entries in its argument. Assume that the benefit function is $[\tau(x)]_i = \tau_i x_i$,
 189 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)
 190 the quotient and the remainder for $\frac{n}{r}$. Let $1 > \bar{\eta} \geq \frac{a\tau_1 - qv}{a\tau_1}$ where $v = a \min \left\{ \sum_{i=1}^q \frac{1}{\tau_i}, \sum_{i=0}^{r-1} \frac{1}{\tau_{n-i}} \right\}$. Then,
 191 $\bar{T} = \left\lceil \frac{n}{r} \right\rceil$ is the smallest number of periods required to achieve perfect fairness.

192 *Proof.* Consider the following allocation: In period t for $1 \leq t \leq \bar{T} - 1$, allocate $\frac{v}{\tau_i}$ for $i = n - (t-1)r - 1, n -$
 193 $(t-1)r - 2, \dots, n - tr$ and 0 to the rest. Note that in each period, we allocate exactly to r stakeholders,
 194 and that the inequality constraint is satisfied due to the definition of v and the ordering of $\tau_1 \geq \tau_2 \dots \geq \tau_n$.
 195 We can ensure that in the first $\bar{T} - 1$ periods, each of the allocated player receives a value of v , and hence
 196 the total benefit of the allocation is rv , giving the $\eta(x(t)) = \frac{a\tau_1 - rv}{a\tau_1} \leq \bar{\eta}$. In the last period, we allocate $\frac{v}{\tau_i}$
 197 for $i = 1, 2, \dots, r$. 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
 198 follows since the benefit to each player is invariant.

199 To prove that the given \bar{T} is the smallest number of periods required to obtain perfect fairness, observe
 200 that, for any fewer number of periods, it is impossible for all players to get a non-zero allocation. ■ □

201 *Remark 15.* Note that, unlike earlier, we do not have tightness in η now.

202 *Remark 16.* We note that the above results, [Theorems 10, 12](#) and [14](#), are agnostic to the choice of ϕ , and
 203 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
 204 necessarily depend upon the choice of ϕ .

205 **Example 17.** Consider [Theorem 14](#) applied to [Example 9](#). We have $a = 1$ and $r = 2$. Again, $\tau(x) =$
 206 $(2x_1, x_2)$, so $\tau_1 = 2$ and $\tau_2 = 1$, and $\tau_1 \geq \tau_2 \geq 0$ holds. $\bar{T} = \left\lceil \frac{n}{r} \right\rceil = \left\lceil \frac{2}{2} \right\rceil = 1$. Indeed, the solution $(\frac{1}{3}, \frac{2}{3})$ is
 207 feasible and fair.

208 4 General Combinatorial \mathcal{X}

209 In many practical areas of interest, \mathcal{X} could be more complicated than the sets presented in [Section 3](#). Let
 210 $\mathcal{X} = \{x^1, x^2, \dots, x^k\}$. We assume that \mathcal{X} is finite, but typically has a large number of elements, given in the
 211 form of a solution to a combinatorial problem. We consider a benefit function where $\tau(x) = \Gamma x$, for some

212 matrix Γ . We thus have $[\tau(x)]_i = \tau_i^\top x$ for $i = 1, \dots, n$.¹ The efficiency constraint here is trivial, as the
 213 inefficient allocations x^j can be removed from \mathcal{X} to retain another (smaller) finite \mathcal{X} . In fact, with linear τ ,
 214 the efficiency constraint is a linear inequality.

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 as²

$$\min_{q,y} \phi(y) \tag{8a}$$

$$\text{s.t. } y_i = \tau_i^\top \left(\sum_{j=1}^k q_j x^j \right), \quad \forall i = 1, \dots, n \tag{8b}$$

$$\sum_{j=1}^k q_j = T \tag{8c}$$

$$q_j \in \mathbb{Z}_{\geq 0}, \quad \forall j = 1, \dots, k. \tag{8d}$$

215 Note that since the theorems of [Section 3](#) do not extend to this case, and since [Example 11](#) shows that a
 216 perfectly fair solution may not even exist, we know neither the value of the fairest feasible solution, nor the
 217 smallest value of T that guarantees such a solution. We now present a two-phase integer program that finds
 218 the smallest value of T that guarantees the fairest feasible solution. In the first phase, we are interested in
 219 finding the fairest feasible solution. This is achieved by solving model (8), replacing (8c) with $\sum_{j=1}^k q_j \geq 1$
 220 to ensure that not only the fairest solution is returned, but that this solution not be $(0, 0, \dots, 0)$. Let ϕ^*
 221 be the value of the solution found in the first phase. In the second phase, we are interested in finding the
 222 smallest T which can accommodate a solution of value ϕ^* . This is achieved by solving model (8), again
 223 replacing (8c) with $\sum_{j=1}^k q_j \geq 1$, adding $\phi(y) = \phi^*$, and changing the objective to $\min_{q,y} \sum_{j=1}^k q_j$.

224 We should note that the value of T^* found by the second phase could be quite high, even for simple
 225 choices of \mathcal{X} . If perfect fairness can be attained in a time horizon of size, say, 10000, and that a discrete time
 226 point represents a day, then this period of 27-odd years would probably not be suitable for most practical
 227 applications, since it is only by reaching the end of the time horizon that optimal fairness is guaranteed. In
 228 practice, then, large time horizons can be problematic. This motivates the need for consistent fairness on a
 229 smaller scale.
 230

231 Some compromises can be made which would mitigate this issue. Manually fixing T to a value not higher
 232 than the expected duration of the problem would ensure a fair distribution of resources, since reaching the
 233 end of the time horizon would be assured. Accepting some degree of unfairness would also decrease the
 234 length of the time horizon and achieve similar results. [Example 18](#) illustrates these two options.

Example 18. 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

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

235 Perfect fairness can be achieved when $T = 1109$, but if we fix $T = 5$, the difference between y_1 and y_2 will
 236 be a mere ~ 0.42 , and if we allow a small difference of 0.01 between y_1 and y_2 , this could be achieved when
 237 $T = 15$.

238 Another compromise would be to choose a fair way of reaching the solution. Given that we know the
 239 value of T^* which grants optimal fairness, we could then identify the *best path* leading to the optimal solution.

¹Matrix Γ can implicitly be seen as an adjacency matrix in order to consider a graph representation of the problem.

²The 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](#).

240 The best path is the one which ensures that unfairness is kept as low as possible in the intermediate time
 241 points, without sacrificing optimality at the end of the time horizon. This would, however, require solving
 242 another linear program, which may prove burdensome if the time horizon were too large.

243 One may think that greedily aiming for the fairest solution at each time point—while considering any
 244 unfairness incurred in previous time points, and enforcing minimal efficiency constraints by ensuring that
 245 all available resources are used at each time point—could lead to optimality at the end of the time horizon.

246 [Theorem 19](#) shows that it is not the case.

247 **Theorem 19.** *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

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

248 where $0 < \epsilon < \frac{1}{2}$. If $T = 2$, it is obvious that the only perfectly fair solution (and, as it happens, the most
 249 efficient too) is achieved by allocating the resource in turn to x_2 and x_3 , over T . Yet, for any reasonable
 250 measure of unfairness, the greedy choice would be to allocate the resource twice to x_1 . Thus, a greedy
 251 approach does not necessarily lead to an optimal solution. ■ □

252 We should note that there is a trade-off between fairness and efficiency: As illustrated in [Example 20](#),
 253 enforcing constraints to increase fairness will generally decrease efficiency, and vice-versa.

Example 20. *Consider a fair allocation 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*

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

254 *By enforcing maximum efficiency, the only feasible solution is $x = (3, 0, 0)$ for $y = (3, \frac{3}{2}, \frac{3}{2})$ and an efficiency
 255 of 6, which is not very fair. In contrast, by enforcing maximum fairness, the only feasible solution is
 256 $x = (0, 1, 1)$ for $y = (1, 1, 1)$ and an efficiency of 3, which is not very efficient (one resource is even wasted
 257 due to the fairness constraints).*

258 We should also note that the value of T may have a noticeable impact on both fairness and efficiency, as
 259 illustrated in [Example 21](#).

Example 21. *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*

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

260 *The smallest T accommodating perfect fairness is $T = 1$ with $x = (1, 0, 1, 0)$ ³ for $y = (1, 1, 1, 1)$ and an
 261 efficiency of 4 (with one wasted resource). However, $T = 2$ also accommodates solutions with perfect fairness,*

³Or any other equivalent solution.

263 *albeit with an increased efficiency. Take, for instance, $x = (3, 0, 0, 0)$ and $x = (0, 0, 0, 3)$ over two time points,*
 264 *for a combined $y = (3, 3, 3, 3)$ and an efficiency of 12. This represents a noticeable increase in efficiency over*
 265 *choosing the initial solution twice to cover the same time horizon.*

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

268 5 Ambulance Location and Relocation

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

272 The definition of this problem is not unified across the literature (Brotcorne et al., 2003, Gendreau et al.,
 273 2001).⁴ Most definitions, however, agree that the population should receive an efficient service, which is
 274 generally translated into some form of coverage. In this paper, we are not solely concerned by efficiency, but
 275 we are further interested in fairness: Ambulances should be allocated such that the same set of people is not
 276 always at a disadvantage with respect to access to a quick service.

277 5.1 Preliminary problem formulation

278 In this manuscript, we adopt the following definition. The region to which ambulances are allocated is
 279 divided into n demand zones, the travel time between each pair being a known quantity. Ambulances can
 280 only be allocated to bases, which are located within a subset of the zones. Each zone (equivalently, the
 281 people living in each zone) are individual stakeholders in this model. A pre-chosen T rounds of decision is
 282 modeled to occur, over which the ambulances should be allocated in a fair and efficient manner. In each
 283 round, a configuration x of how m ambulances are placed is chosen. Next, we define the benefit function
 284 for the zones. In this model, each zone i , based on its population, has a demand of ζ_i ambulance vehicles.
 285 A zone i is said to be *covered* in configuration x if there are at least ζ_i ambulances located in bases from
 286 which zone i could be reached in a time less than a chosen time limit called the *response threshold time*.
 287 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
 288 a nonlinear benefit function in this case, as opposed to the simpler cases in Section 3 to both reflect the fact
 289 that more number of ambulances are essential to sufficiently serve certain regions and also to demonstrate
 290 the robustness of our methods to even certain classes of nonlinear (piecewise linear) benefit functions. The
 291 nonlinearity, as shown below, can be modeled using auxiliary variables and linear functions, to conform to
 292 the form of (8). Furthermore, efficiency constraints analogous to (7c) are enforced by stating that at least a
 293 fraction f of the zones should be covered in each allocation. The unfairness metric ϕ used is the difference
 294 between the two zones which are most often and least often covered, over T .

295 With the above, the problem is a standard T-PFA problem in Definition 8. However, the ambulance
 296 allocation problem involves additional constraints on allocations between two consecutive periods. Decision-
 297 makers prefer policies that ensure that not too many ambulances have to shift bases on a daily basis. Thus,
 298 between two consecutive allocations, we would ideally like to have not more than a fixed number r of
 299 ambulances to shift bases. We refer to this constraint as the *transition constraint*.

Now, the problem can be formally cast as follows:

$$\min_{y, v, x, \tau, \phi} \phi(y) \tag{9a}$$

$$\text{s.t. } x_i(t) = 0 \quad \forall i \notin \mathcal{B}; \forall t \in \mathcal{T} \tag{9b}$$

$$\sum_{i=1}^n x_i(t) \leq m \quad \forall t \in \mathcal{T} \tag{9c}$$

⁴This 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.

$$v_i(t) = \sum_{j=1}^n a_{ji} x_j(t) \quad \forall i = 1, \dots, n; t \in \mathcal{T} \quad (9d)$$

$$v_i(t) \leq (\zeta_i - 1) + m \tau_i(t) \quad \forall i = 1, \dots, n; t \in \mathcal{T} \quad (9e)$$

$$v_i(t) \geq \zeta_i - m(1 - \tau_i(t)) \quad \forall i = 1, \dots, n; t \in \mathcal{T} \quad (9f)$$

$$y_i = \frac{1}{T} \sum_{t=1}^T \tau_i(t) \quad \forall i = 1, \dots, n \quad (9g)$$

$$\sum_{i=1}^n \tau_i(t) \geq fn \quad \forall t \in \mathcal{T} \quad (9h)$$

$$x_i(t) \in \mathbb{Z}_{\geq 0} \quad \forall i = 1, \dots, n; t \in \mathcal{T} \quad (9i)$$

$$\tau_i(t) \in \{0, 1\} \quad \forall i = 1, \dots, n; t \in \mathcal{T} \quad (9j)$$

$$\sum_{i=1}^n |x_i(t+1) - x_i(t)| \leq 2r \forall t \in 1, \dots, T-1 \quad (9k)$$

300

301 Here, $\mathcal{T} = \{1, \dots, T\}$, $x_i(t)$ is the number of ambulances allotted to zone i at time t . Constraints (9b) ensure
 302 that allocation occurs only if i is an ambulance base. Here, \mathcal{B} is the set of ambulance bases. Constraints (9c)
 303 limit the number of ambulances available for allocation in each round. Moreover, a_{ij} is a binary parameter
 304 which is 1 if an ambulance can go from i to j within the response threshold time and 0 otherwise. Through
 305 constraints (9d), $v_i(t)$ counts the number of ambulances that can reach zone i within the response threshold
 306 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
 307 otherwise. This is accomplished by constraints (9e) and (9f). Constraints (9h) take care of efficiency and
 308 ensure only allocations that cover at least a fraction f of all zones are to be considered. Finally, constraints
 309 (9k) are the transition constraints, which can easily be reformulated through linear inequalities. It ensures
 310 that not too many ambulances shift bases between consecutive time periods.

311 The problem in (9) is a mixed-integer linear program (MILP). However, the problem is symmetric with
 312 respect to certain permutations of the variables. Symmetry makes it particularly hard for modern branch-
 313 and-bound-based solvers (Margot, 2010). Even if one relaxes the transition constraints (9k), the symmetry
 314 exists. One can observe that, relaxing the transition constraints (9k) in (9), we have a T-PFA problem. We
 315 call this relaxed problem defined by (9a) to (9j) as the Ambulance-without-transition constraints (AWT)
 316 problem.

317 5.2 A branch-and-price reformulation of the AWT

318 Since the AWT in (9a) to (9j) is a T-PFA problem, it can be readily written in the form shown in (8) and
 319 hence can be solved using branch-and-price. First, we note that without the constraint in (9k), any feasible
 320 solution or optimal solution remains feasible or optimal after permutations to t . Thus, the following version
 321 of the problem could be used to solve the relaxed problem without the symmetry. Here we only *count* the
 322 number of times each configuration might be used over T periods.

The master problem (MP).

$$\min \quad g - h \quad (10a)$$

$$\text{s.t.} \quad g \geq y_i \quad (\alpha_i) \quad \forall i = 1, \dots, n \quad (10b)$$

$$y_i \geq h \quad (\beta_i) \quad \forall i = 1, \dots, n \quad (10c)$$

$$y_i = \frac{1}{T} \sum_{j=1}^k \tau_{ij} q_j \quad (\lambda_i) \quad \forall i = 1, \dots, n \quad (10d)$$

$$\sum_{j=1}^k q_j = T \quad (\mu) \quad (10e)$$

$$q_j \geq 0 \quad \forall j = 1, \dots, k \quad (10f)$$

$$q_j \in \mathbb{Z} \quad \forall j = 1, \dots, k. \quad (10g)$$

323

324 where q_j counts the number of times the configuration defined by x^j is used. The benefit obtained by
 325 stakeholder i due to the configuration x^j is τ_{ij} . y_i is the average benefit that stakeholder i enjoys through
 326 the time horizon of planning. Note that, once we solve (10), an equivalent solution to the AWT could be
 327 obtained by arbitrarily considering the allocations x^j for q_j times in $x(1), \dots, x(T)$ of (9). Similarly given a
 328 solution to AWT, one could immediately identify a corresponding solution to (10).

329

330 However, since the number of configurations are typically exponentially large and we only might use a
 331 handful of them in a solution, we could resort to a branch-and-price approach where (10) only contains a
 subset of the configurations.

Referring to the continuous relaxation of MP as CMP, we write the dual of CMP below.

$$\max T\mu \quad (11a)$$

$$\text{s.t.} \quad \sum_{i=1}^n \alpha_i = 1 \quad (g) \quad (11b)$$

$$\sum_{i=1}^n \beta_i = 1 \quad (h) \quad (11c)$$

$$\lambda_i - \alpha_i + \beta_i = 0 \quad (y_i) \quad \forall i = 1, \dots, n \quad (11d)$$

$$\mu \leq \frac{1}{T} \sum_{i=1}^n \tau_{ij} \lambda_i \quad (q_j) \quad \forall j = 1, \dots, k \quad (11e)$$

$$\alpha_i, \beta_i \geq 0 \quad \forall i = 1, \dots, n. \quad (11f)$$

332

333 Considering only a subset of columns in the CMP is equivalent to considering only a subset of constraints
 334 in (11e). Given some dual optimal solution $(\alpha^*, \beta^*, \lambda^*, \mu^*)$ to the dual of the restricted CMP, one can find
 the most violated constraint in (11e) and include the corresponding column in the restricted CMP.

The pricing problem.

$$\min \sum_{i=1}^n \tau_i \lambda_i^* \quad (12a)$$

$$\text{s.t.} \quad (\tau_i = 1) \iff (a_i^\top x \geq \zeta_i), \quad \forall i = 1, \dots, n \quad (12b)$$

$$(\tau_i = 0) \iff (a_i^\top x \leq \zeta_i - 1), \quad \forall i = 1, \dots, n \quad (12c)$$

$$\sum_{i=1}^n x_i \leq m \quad (12d)$$

$$\sum_{i=1}^n \tau_i \geq fn \quad (12e)$$

$$\tau \in \{0, 1\}^n \quad (12f)$$

$$x \in \mathbb{Z}_{\geq 0}^n. \quad (12g)$$

335

The minimal efficiency constraints are embedded in the pricing problem. The time horizon is fixed in (10e).

336 Constraints (12b) and (12c) help compute the benefits to each zone, τ and can clearly be rewritten
 337 with integer variables and linear constraints. No more than m ambulances may be used (12d), and the
 338 configuration must cover at least a fraction f of the zones (12e).

339 The reformulation (10) of the AWT in (9a) - (9j), can now be solved very efficiently using branch-and-
 340 price. However, a solution thus obtained might violate (9k). Further, given the solution in the space of
 341 variables in (10) it is not immediate if one can easily verify whether the constraint (9k) is or is not satisfied.
 342 So, given a feasible point for (10), we define the configuration graph as follows, and then show that the point
 343 satisfies (9k) if and only if the configuration graph has a Hamiltonian path.

344 5.3 Checking (9k)

345 **Definition 22** (Configuration graph). *Given a solution \bar{q} to (10), define $\mathcal{I} = \{j : \bar{q}_j \geq 1\}$. The con-*
 346 *figuration 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*
 347 *$E = \{((j_1, j'_1), (j_2, j'_2)) : \|\bar{x}^{j_1} - \bar{x}^{j_2}\|_1 \leq 2r\}$ where \bar{x}^{j_1} and \bar{x}^{j_2} are the configurations corresponding to*
 348 *the variables \bar{q}_{j_1} and \bar{q}_{j_2} .*

349 **Theorem 23.** *Given a point \bar{q} that is feasible to (10), there exists a corresponding \bar{x} that is feasible to (9)*
 350 *if the CG defined by \bar{q} contains a Hamiltonian path. Conversely, if there is a feasible solution to (9) which*
 351 *uses only the configurations with indices in $\mathcal{I} = \{j : \bar{q}_j \geq 1\}$, then the corresponding CG has a Hamiltonian*
 352 *path.*

353 *Proof.* Observe that G has exactly T vertices, since if a configuration is found more than once in \bar{q} , it is split
 354 into as many distinct vertices in V , which are distinguished by the second element of the tuple. An edge
 355 $((u_1, u_2), (v_1, v_2))$ in E indicates that movement between configurations indexed by u_1 and v_1 does not violate
 356 the transition constraints. A Hamiltonian path is a path that visits each vertex exactly once. As such, any
 357 Hamiltonian path in G proves that a feasible sequence of transitions which does not violate the transition
 358 constraints exists between the configurations in \bar{q} . Given a Hamiltonian path $(v_1, w_1), (v_2, w_2), \dots, (v_T, w_T)$
 359 in G , a feasible sequence $x(1), x(2), \dots, x(T)$ for (9) would be $\bar{x}^{v_1}, \bar{x}^{v_2}, \dots, \bar{x}^{v_T}$.

360 Conversely, given a feasible sequence $x^{j_1}, x^{j_2}, \dots, x^{j_T}$ for (9), a Hamiltonian path can be constructed for
 361 G as $(j_\alpha, \beta) : \alpha \in \{1, \dots, T\}, \beta = 1 + \max_{\beta'}((j_{\alpha'}, \beta') : \alpha' \in \{1, \dots, \alpha - 1\} \wedge x^{j_\alpha} = x^{j_{\alpha'}})$. ■ □

362 Following Theorem 23, one can construct the CG with exactly T vertices and can check if the solution
 363 to (10) is feasible to (9). If yes, we are done. If not, the way we can proceed is detailed in the rest of this
 364 section.

365 5.3.1 Cutting planes and binarization

366 The most natural way to eliminate a solution that does not satisfy a constraint is by means of a cutting
 367 plane. This is a common practice in the computational MILP research, for example, the subtour elimination
 368 inequalities in a travelling salesman problem.

369 In cases where a more sophisticated cutting plane is not available, but every feasible solution is determined
 370 by a binary vector, no-good cuts could be used to eliminate infeasible solutions one by one (D'Ambrosio
 371 et al., 2010). However, the variables x in (10) are general integer variables. It is possible that a point x that
 372 violates the transition constraint could lie strictly in the convex hull of solutions that satisfy the transition
 373 constraint. Hence it could be impossible to separate x using a cutting plane. Example 24 demonstrates the
 374 above phenomenon.

375 **Example 24.** *Consider the case where (10) and (12) has an optimal solution given by $q = (q_1, q_2, q_3, q_4) =$*
 376 *$(1, 1, 1, 1)$ and the allocations x_1, x_2, x_3, x_4 (in the context of (9)) corresponding to q_1, q_2, q_3, q_4 are $(3, 3, 3, 3),$*
 377 *$(4, 2, 3, 3), (2, 4, 3, 3), (3, 3, 2, 4)$ respectively. Let $r = 1$. Then, the CG corresponding to this solution is a tree*
 378 *as shown in Figure 1 and hence does not have a Hamiltonian path. On the other hand, for the choices $q^i,$*
 379 *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*
 380 *and hence have a Hamiltonian path, trivially.*

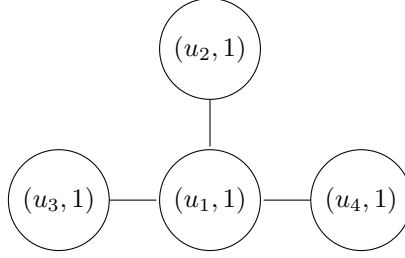


Figure 1: Configuration graph for the solution in [Example 24](#).

381 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
 382 solutions satisfies the transition constraint, q does not. Thus no valid cutting plane can cut only the infeasible
 383 point.

384 While [Example 24](#) shows that an infeasible solution cannot always be separated using cutting planes, it
 385 could still be possible that there is an extended formulation where the infeasible point could be separated.

386 **Naive binarization.** A natural choice of an extended formulation comes from binarizing each integer
 387 variable q_j in (10). This is possible because each q_j is bounded above by T and below by 0. Thus one can
 388 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 power of
 389 2 lesser than or equal to T . If each b_j^ℓ is binary, any integer between 0 and T can be represented as above.
 390 Having written the above binarization scheme, one can separate any solution q by adding a no-good cut on
 391 the corresponding binary variables. A no-good cut is a linear inequality that separates a single vertex of
 392 the 0 – 1 hypercube without cutting off the rest ([Balas and Jeroslow, 1972](#), [D’Ambrosio et al., 2010](#)). An
 393 example is shown in [Example 25](#).

Example 25. Consider the problem in [Example 24](#). Binarization to separate the solution $q = (1, 1, 1, 1)$ can
 be done by adding the following constraints to (10).

$$q_j = b_j^1 + 2b_j^2 + 4b_j^4 \quad \forall j = 1, \dots, 4 \quad (13a)$$

$$\sum_{j=1}^4 (b_j^1 + (1 - b_j^2) + (1 - b_j^4)) \leq 11 \quad (13b)$$

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

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

400 **Strengthened binarization.** When the CG corresponding to q is not just lacking a Hamiltonian path,
 401 but is a disconnected graph (a stronger property holds), one could add a stronger cut, which could potentially
 402 cut off multiple infeasible solutions. In this procedure, we add a binary variable b_j for each q_j such that
 403 $b_j = 1$ if and only if $q_j \geq 1$. Now, observe that if a CG corresponding to the solution \bar{q} is disconnected, then
 404 it will be disconnected for all q whose non-zero components coincide with the non-zero components of \bar{q} . i.e.,
 405 no solution with the same support as that of \bar{q} could satisfy the transition constraints. Thus one could add
 406 a no-good cut on these b_j binary variables, which cuts off all the solutions with the same support as \bar{q} .

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

410 5.3.2 Three-way branching

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

415 In this method, as soon as a solution \bar{q} satisfying all the integrality constraints but violating the transition
 416 constraint is found, the following actions are performed. First, if the lower bound and the upper bound for
 417 every component of \bar{q} match, then we are at a leaf that can be discarded as infeasible. If not, find a component
 418 j (in our case, a configuration j) such that \bar{q}_j is strictly different from at least one of the bounds. Now, we
 419 do a three-way branching on the variable q_j where the new constraints in each of the three branches are (i)
 420 $q_j \leq \bar{q}_j - 1$ (ii) $q_j = \bar{q}_j$ (iii) $q_j \geq \bar{q}_j + 1$.

421 Finite termination follows from the fact that \bar{q}_j is infeasible for branches (i) and (iii) and hence will never
 422 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
 423 have one more fixed variable and this variable will never be branched on again.

424 Three-way branching is used as opposed to regular two-way branching but on integer variables because
 425 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
 426 cut off other feasible solutions with $q_j = \bar{q}_j$ but the solution differing in components other than j . Branching
 427 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,
 428 nor have we eliminated the infeasible solution. Thus, the LP optimum in the second branch is going to be \bar{q}
 429 again, and cycling ensues.

430 5.3.3 Constraint Programming

431 The final alternative we can use to enforce transition constraints (9k) is constraint programming (CP).
 432 Namely, CP is a programming paradigm for solving combinatorial problems. A CP model is defined by
 433 a set of variables, each of which is allowed values from a finite set, its *domain*. The relationship between
 434 these variables is determined by the constraints of the problem. These succinct constraints can encapsulate
 435 complex combinatorial structures, such as the packing of items into bins, for example. A solver then solves
 436 the problem by enforcing consistency between the variables and the constraints, and using branching and
 437 backtracking techniques to explore the solution space. Before defining the CP model to enforce constraints
 438 (9k), we define the compact configuration graph, and explain its relationship with the CG.

439 **Definition 26** (Compact configuration graph). *Given a feasible solution q^* to the continuous relaxation of*
 440 *(10) (CMP), the compact configuration graph (CCG) is defined as $G = (V, E)$, where $V = \mathcal{A} := \{j : q_j^* > 0\}$*
 441 *and $E = \{(v, w) : \|q_v^* - q_w^*\|_1 \leq 2r\}$.*

442 **Definition 27** (Walk). *A walk in an undirected graph $G = (V, E)$ is a finite sequence of vertices v_1, \dots, v_k ,*
 443 *not necessarily distinct, such that for each i , $(v_i, v_{i+1}) \in E$.*

444 **Theorem 28.** *Every walk of length T in a CCG constructed from a solution q^* to CMP corresponds to a*
 445 *feasible solution of (9).*

446 *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 .
 447 Since we allow revisiting vertices, it is possible that $v_j = v_{j'}$ for some $j \neq j'$.

448 Now define \tilde{q} component-wise where \tilde{q}_j corresponds to the number of times the vertex j is visited in the
 449 walk W . Since the walk has a length T , trivially $\sum_j \tilde{q}_j = T$, satisfying (10e). Then, \tilde{y}, \tilde{q} can be defined so that
 450 we have a feasible solution to (10). Hence, if we now show that the CG defined by the nonzero components
 451 of \tilde{q} has a Hamiltonian path, then the corresponding \tilde{x} will be feasible for (9) due to Theorem 23.

452 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 $n(v_j)$
 453 records the number of times v_j has appeared in the path P earlier so far. We note that each term in the path
 454 P is indeed a vertex of the CG (which has T vertices) and that they are all visited exactly once, implying
 455 that P is the required Hamiltonian path. ■ □

456 The general mechanism involving the CP component is provided in [Algorithm 1](#), which is the entire
457 algorithm we test in [Section 6](#). The CP component checks whether the solution returned by the CMP
458 can be made valid in some way, i.e., if there exist solutions using only the configurations in \mathcal{A} (i.e., the
459 configurations that appear in the optimal CMP solution, see [Definition 26](#)) with integer values, such that
460 they do not violate the transition constraints. By providing these feasible solutions, it provides an upper
461 bound to (9). As soon as a set of configurations is given to CP, a cut is added to the CMP, which eliminates
462 all solutions to CMP which only consist of the configurations provided to CP. Namely, if \mathcal{A} indices the
463 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
464 must be outside the set indexed by \mathcal{A} .

465 Let k' be the cardinality of \mathcal{A} . A CCG is associated with solution q — this CCG forms the basis for a
466 deterministic finite automaton. This automaton A is defined by a tuple $(Q, \Sigma, \delta, q_0, F)$ of states Q , alphabet
467 Σ , transition function $\delta : Q \times \Sigma \rightarrow Q$, initial state $q_0 \in Q$, and final states $F \subseteq Q$, with $Q = \{0, \dots, k'\}$,
468 $\Sigma = \{1, \dots, k'\}$, $\delta = \{(u, v) \rightarrow v : (u, v) \in E\} \cup \{(0, u) \rightarrow u : u \in F\}$, $q_0 = 0$, and $F = \{1, \dots, k'\}$. In other
469 words, this automaton accepts any valid sequence of k' configurations, with dummy configuration $q_0 = 0$
470 being the initial state. Let LB be the lower bound given by the CMP, and UB be an upper bound. Finally,
471 let Ω be the collection of all index sets \mathcal{A}' for which the CP model has been previously solved. There are
472 T decision variables z with domains $\{0, \dots, k' - 1\}$, with z_t indicating which CMP configuration from \mathcal{A} is
473 used at time point t . These variables are constrained by

$$\min \phi \tag{14a}$$

$$\phi = \max(c) - \min(c) \tag{14b}$$

$$LB \leq \phi \leq UB - 1 \tag{14c}$$

$$\text{cost_regular}(z, A, \tau_{i*}, c_i) \quad i = 1, \dots, n \tag{14d}$$

$$\text{at_most}(T - 1, z, \omega) \quad \forall \omega \subset \mathcal{A} : \exists \mathcal{A}' \in \Omega, \omega \subseteq \mathcal{A}' \tag{14e}$$

$$z_t \in \mathcal{A} \quad t = 1, \dots, T \tag{14f}$$

$$c_i \in \mathbb{Z}_{\geq 0} \quad i = 1, \dots, n \tag{14g}$$

474 The objective (14a) is to minimize the unfairness, i.e., the difference between the zone which is covered
475 the most, and that which is covered the least (14b). Any feasible solution should be strictly better than
476 the upper bound, and search by the CP solver can be interrupted as soon as a feasible solution is found
477 whose objective value coincides with LB (14c). By [Theorem 28](#), any sequence of configurations indexed
478 by \mathcal{A} and of length T must correspond to a feasible solution of (9). This requirement is enforced by the
479 **cost-regular** ([Demasse et al., 2006](#)) constraints (14d): A **cost-regular** constraint holds for zone i if
480 z forms a word recognized by automaton A and if variable c_i is equal to the coverage of zone i over T :
481 $c_i = \sum_{t=1}^T \tau_{i, z_t}$. Note that the CP model does not require that all configurations indexed by \mathcal{A} be used at
482 least once: Since it checks all possible subsets of \mathcal{A} , the cut added to CMP does not remove any unexplored
483 part of the search space. Finally, all subsets $\omega \subset \mathcal{A}$ previously explored, i.e. being as well a subset of some
484 index set \mathcal{A}' in Ω , are considered by the model using the **at_most** constraint (14e): At most $T - 1$ variables
485 in z can take on values in ω . The motivation for this is illustrated in [Example 29](#).

487 **Example 29.** Assume that $\mathcal{A} = \{j', j'', j'''\}$. If (14) had previously solved $\mathcal{A}' = \{j', j''\}$, this solution space
488 would be explored again since the configurations in \mathcal{A} are not constrained to be used at least once. Adding
489 constraint **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 (14b) and (14d) with the equivalent but less efficient⁵ reformulation

$$\phi = \max_{i=1}^n \sum_{t=1}^T \tau_{i, z_t} - \min_{i=1}^n \sum_{t=1}^T \tau_{i, z_t}$$

⁵In practice, the CP component of [Algorithm 1](#) remains very fast, so this loss in efficiency is negligible.

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 (9).

```
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 (10) 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 (14) on the graph  $G$  with appropriate values of  $\tau, c.$ 
17:   if (14) 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
```

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 (10) 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 (9) 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,

505 we use real data from the city of Utrecht, Netherlands, as well as synthetically-generated instances of varying
506 sizes. Moreover, we consider a predetermined time horizon, i.e., a fixed T of size 30, with the understanding
507 that a decision maker could reasonably make plans on a monthly basis.

508 **Utrecht instance.** We use the model defined in Section 5 to determine ambulance allocation in the city
509 of Utrecht, Netherlands, using a dataset provided by the RIVM.⁶ The Utrecht instance contains 217 zones,
510 18 of which are bases. Since calls should be reached within 15 minutes, we consider that a zone is connected
511 to another if an ambulance can travel between the two zones in under 15 minutes. This makes the graph
512 about 28.4% dense.

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

517 **Synthetic instances.** Reproducing the ratio of bases and edges to zones, we also generated synthetic
518 instances⁷ of sizes 50, 100, 200, and 400 zones, using the Matérn cluster point process (Matérn, 1960), which
519 helps in generating a distribution of zones mimicking a realistic urban setting. The number of ambulances
520 for the instances is chosen to be just enough to ensure feasibility. The results are averaged over five instances
521 of each size.

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

525 **Results.** It is easy to see that there are two parameters that are likely to influence the difficulty of the
526 solving process: the size of the instance (number of zones) and the transition constraints, i.e., the flexibility
527 we allow for relocating ambulances between zones on a daily basis. We would typically expect the size of the
528 instance to be a significant factor, but this is likely to be mitigated by the column generation approach, which
529 tends to scale well. The effects that the transition constraints will have on the solving process, however, are
530 less clear. In order to assess such an effect, for each instance, we vary the number of ambulances that can
531 shift bases on consecutive days (r in constraints (9k)) from $0.1m$ to m in increments of $0.1m$, where m is
532 the total number of ambulances in the instance. We call this ratio maximum transition, MT . For each of
533 these cases, we record the average of the time taken to solve these instances to optimality (capped at 1,200
534 CPU seconds). We record the final relative gap left for the instance, which is defined by $\frac{|UB_f - LB_f|}{UB_f}$, where
535 LB_f and UB_f represent the values of the best lower and upper bounds at the time limit, respectively. We
536 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
537 value of the initial lower bound (without any cuts), and UB_i represents the value of the initial, trivial upper
538 bound. We present the final relative gaps associated with different classes of instances in Figure 2.

539 One can immediately observe from Figure 2 that the transition constraints are affecting the difficulty of
540 the instances for $MT \leq 0.5$, while everything can be solved to optimality within the time limit for $MT \geq 0.5$
541 independently of the number of zones. The larger MT the smaller the final gap, i.e., when there is more
542 freedom in moving the ambulances around on a daily basis, Algorithm 1 becomes very effective in computing
543 optimal solutions.

544 We observe that for varying values of MT , the Utrecht instance has a higher relative gap than synthetic
545 instances of comparable and even larger size. This discrepancy is the result of the distribution of the
546 population densities across the zones. In the synthetic instances, the population densities are randomly

⁶National Institute for Public Health and the Environment of the Netherlands.

⁷These instances can be found at <https://github.com/PhilippeOlivier/ambulances>.

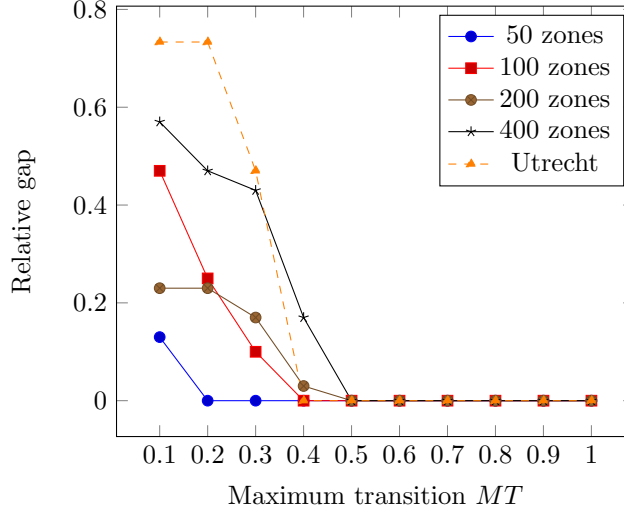


Figure 2: Final relative gaps with respect to the maximum transition MT . Gaps in the synthetic instances are averaged over 5 instances. Time limit of 1,200 CPU seconds.

distributed among the zones. The Utrecht instance, in contrast, exhibits several distinct clusters of varying sizes and population density distributions.⁸

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 1, 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”),⁹ and the average time for the instances solved to optimality (“time”).

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 1: 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.

The results in Table 1 show that for $MT > 0.7$ constraints (9k) 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 $MT = 0.7$.

⁸Constructing 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.

⁹The entry “#solved” takes an integer value between 0 and 5 for synthetic instances and 0/1 for the Utrecht instance.

556 The average time for the instances solved to optimality is often low, thus indicating that if we manage to
 557 prove optimality by completely closing the gap, then this is achieved quickly, generally with a few calls to
 558 CONSTRAINTPROGRAMMING. In contrast, if we do not succeed in closing the gap early, then it is unlikely to
 559 be closed at all in the end. In the cases where optimality is not proven, improvements are still made during the
 560 solving process mostly because of improvements on the upper bound value as CONSTRAINTPROGRAMMING
 561 tests more and more combinations of configurations. This suggests a few things. First, with a high enough
 562 MT , the initial lower bound is generally optimal, and feasible solutions whose objective values coincide with
 563 this bound can readily be found. Second, when the MT value is low, either feasible solutions whose objective
 564 values coincide with the initial lower bound are very rare, or they are nonexistent. Third, the computational
 565 power of Algorithm 1 is associated with the effectiveness in searching for primal solutions (upper bound
 566 improvement) versus the progress on the dual side (lower bound improvement), which appears to be very
 567 difficult.

568 Finally, we are also interested in identifying the time spent in the column generation part as opposed to
 569 the CONSTRAINTPROGRAMMING routine. To this end, we track the number of calls made to the function
 570 CONSTRAINTPROGRAMMING and also measure the total time spent in calls to the function. The number of
 571 generated columns and the number of calls to CONSTRAINTPROGRAMMING are detailed in Table 2.

MT	50 zones		100 zones		200 zones		400 zones		Utrecht	
	#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 2: Number of columns and calls to CONSTRAINTPROGRAMMING associated with the maximum transition (MT) constraints of the various instance sizes.

572 The results in Table 2 show that when enough freedom in the day-to-day movement of ambulances is
 573 allowed, optimality can generally be proven immediately, with just one call to the CONSTRAINTPROGRAMMING
 574 routine. We also notice that the ratio of columns to the number of calls to the routine increases
 575 with the instance size. Such a ratio for the Utrecht instance is disproportionately high, due again to the
 576 distribution of the population densities.

577 Finally, Table 3 shows the average amount of time spent generating columns (“CG”), that spent enforcing
 578 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 3: CPU times in seconds spent in column generation and CONSTRAINTPROGRAMMING, with $MT = 0.1$.

579 The results in Table 3 show that generating columns for larger instances (more zones) takes naturally
 580 more time, which is why this effort increases with n . The enforcement of the transition constraints, however,
 581

582 requires a more stable effort, as the input of CONSTRAINTPROGRAMMING is always around three to ten
583 variables, regardless of the size of the instance.

584 7 Conclusion

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

592 Acknowledgements

593 The authors would like to thank the National Institute for Public Health and the Environment of the
594 Netherlands (RIVM) for access to the data of their ambulance service.

595 References

- 596 Egon Balas and Robert Jeroslow. Canonical cuts on the unit hypercube. *SIAM Journal on Applied Mathe-*
597 *matics*, 23(1):61–69, 1972.
- 598 Evripidis Bampis, Bruno Escoffier, and Sasa Mladenovic. Fair resource allocation over time. In *Proceedings*
599 *of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, AAMAS '18, page
600 766–773, Richland, SC, 2018. International Foundation for Autonomous Agents and Multiagent Systems.
- 601 Luce Brotcorne, Gilbert Laporte, and Frédéric Semet. Ambulance location and relocation models.
602 *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>.
- 603
604
- 605 Oscar R. Burt and Curtis C. Harris. Letter to the editor—apportionment of the u.s. house of representatives:
606 A minimum range, integer solution, allocation problem. *Operations Research*, 11(4):648–652, 1963. doi:
607 10.1287/opre.11.4.648. URL <https://doi.org/10.1287/opre.11.4.648>.
- 608 Sophie Demassez, Gilles Pesant, and Louis-Martin Rousseau. A cost-regular based hybrid column generation
609 approach. *Constraints*, 11(4):315–333, Dec 2006. ISSN 1572-9354. doi: 10.1007/s10601-006-9003-7. URL
610 <https://doi.org/10.1007/s10601-006-9003-7>.
- 611 Claudia D’Ambrosio, Antonio Frangioni, Leo Liberti, and Andrea Lodi. On interval-subgradient and no-good
612 cuts. *Operations Research Letters*, 38(5):341–345, 2010.
- 613 Ohad Eisenhandler and Michal Tzur. The humanitarian pickup and distribution problem. *Operations*
614 *Research*, 67(1):10–32, 2019. doi: 10.1287/opre.2018.1751. URL <https://doi.org/10.1287/opre.2018.1751>.
- 615
- 616 Gerald Gamrath, Daniel Anderson, Ksenia Bestuzheva, Wei-Kun Chen, Leon Eifler, Maxime Gasse,
617 Patrick Gemander, Ambros Gleixner, Leona Gottwald, Katrin Halbig, Gregor Hendel, Christopher Ho-
618 jny, Thorsten Koch, Pierre Le Bodic, Stephen J. Maher, Frederic Matter, Matthias Miltenberger, Erik
619 Mühmer, Benjamin Müller, Marc E. Pfetsch, Franziska Schösser, Felipe Serrano, Yuji Shinano, Christine
620 Tawfik, Stefan Vigerske, Fabian Wegscheider, Dieter Weninger, and Jakob Witzig. The SCIP Optimization

621 Suite 7.0. ZIB-Report 20-10, Zuse Institute Berlin, March 2020a. URL <http://nbn-resolving.de/urn:nbn:de:0297-zib-78023>.

622

623 Gerald Gamrath, Daniel Anderson, Ksenia Bestuzheva, Wei-Kun Chen, Leon Eifler, Maxime Gasse,
624 Patrick Gemander, Ambros Gleixner, Leona Gottwald, Katrin Halbig, Gregor Hendel, Christopher Hojny,
625 Thorsten Koch, Pierre Le Bodic, Stephen J. Maher, Frederic Matter, Matthias Miltenberger, Erik Mühmer,
626 Benjamin Müller, Marc E. Pfetsch, Franziska Schlösser, Felipe Serrano, Yuji Shinano, Christine Tawfik,
627 Stefan Vigerske, Fabian Wegscheider, Dieter Weninger, and Jakob Witzig. The SCIP Optimization Suite
628 7.0. Technical report, Optimization Online, March 2020b. URL http://www.optimization-online.org/DB_HTML/2020/03/7705.html.

629

630 Michel Gendreau, Gilbert Laporte, and Frédéric Semet. A dynamic model and parallel tabu search heuristic
631 for real-time ambulance relocation. *Parallel Computing*, 27(12):1641–1653, 2001. ISSN 0167-8191.
632 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.

633

634 Jessica L. Heier Stamm, Nicoleta Serban, Julie Swann, and Pascale Wortley. Quantifying and explaining
635 accessibility with application to the 2009 h1n1 vaccination campaign. *Health Care Management Science*,
636 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>.

637

638 Stephen Jacobsen. On marginal allocation in single constraint min-max problems. *Management Science*, 17
639 (11):780–783, 1971. doi: 10.1287/mnsc.17.11.780. URL <https://doi.org/10.1287/mnsc.17.11.780>.

640 Naoki Katoh, Toshihide Ibaraki, and H. Mine. An algorithm for the equipollent resource allocation problem.
641 *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>.

642

643 Bernard O. Koopman. The optimum distribution of effort. *Journal of the Operations Research Society of*
644 *America*, 1(2):52–63, 1953. doi: 10.1287/opre.1.2.52. URL <https://doi.org/10.1287/opre.1.2.52>.

645 Hanan Luss. *Equitable Resource Allocation: Models, Algorithms and Applications*. Information and Commu-
646 nication Technology Series,. Wiley, 2012. ISBN 9781118449219. URL https://books.google.ca/books?id=Z2_3oWjnASkC.

647

648 François Margot. Symmetry in integer linear programming. In *50 Years of Integer Programming 1958-2008*,
649 pages 647–686. Springer, 2010.

650 Bertil Matérn. Spatial variation – stochastic models and their applications to some problems in forest survey
651 sampling investigations. *Report of the Forest Research Institute of Sweden*, 49(5):1–144, 1960.

652 Włodzimierz Ogryczak, Hanan Luss, Michał Pióro, Dritan Nace, and Artur Tomaszewski. Fair optimization
653 and networks: A survey. *Journal of Applied Mathematics*, 2014:612018, Sep 2014. ISSN 1110-757X. doi:
654 10.1155/2014/612018. URL <https://doi.org/10.1155/2014/612018>.

655 Ogryczak, Włodzimierz. Fair optimization – methodological foundations of fairness in network resource
656 allocation. In *2014 IEEE 38th International Computer Software and Applications Conference Workshops*,
657 pages 43–48, 2014.

658 Evan L. Porteus and Jonathan S. Yormark. More on min-max allocation. *Management Science*, 18(9):
659 502–507, 1972. doi: 10.1287/mnsc.18.9.502. URL <https://doi.org/10.1287/mnsc.18.9.502>.

660 Naoki Katoh Toshihide Ibaraki. *Resource Allocation Problems: Algorithmic Approaches*, page 1. Foundations
661 of Computing. The MIT Press, 1988. ISBN 0262090279,9780262090278.

662 Zeev Zeitlin. Minimization of maximum absolute deviation in integers. *Discrete Applied Mathematics*, 3
663 (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>.

664