

---

# OPTIMAL RESOURCE AND DEMAND REDISTRIBUTION FOR HEALTHCARE SYSTEMS UNDER STRESS FROM COVID-19

---

**Felix Parker**

Department of Civil and Systems Engineering  
Johns Hopkins University, Baltimore, MD 21218  
fparker9@jhu.edu

**Hamilton Sawczuk**

Department of Computer Science  
Johns Hopkins University, Baltimore, MD 21218  
hsawczu1@jhu.edu

**Fardin Ganjkanloo**

Department of Civil and Systems Engineering  
Johns Hopkins University, Baltimore, MD 21218  
fganjkh1@jhu.edu

**Farzin Ahmadi**

Department of Civil and Systems Engineering  
Johns Hopkins University, Baltimore, MD 21218  
fahmadi1@jhu.edu

**Kimia Ghobadi**

Department of Civil and Systems Engineering  
Johns Hopkins University, Baltimore, MD 21218  
kimia@jhu.edu

## ABSTRACT

**Problem:** When facing an extreme stressor, healthcare systems typically respond reactively by creating surge capacity at facilities that are at or approaching their baseline capacity. The amount of total surge capacity required can be reduced, however, by redistributing demand and critical resources between facilities to take advantage of underutilized capacity. We study the problem of finding optimal patient and resource transfers to minimize the required surge capacity and resource shortage during a period of heightened demand.

**Academic/Practical relevance:** The COVID-19 pandemic has added a substantial burden to many healthcare systems and has threatened to overwhelm hospitals. Data shows that this additional load was unevenly distributed between hospitals, requiring some to create surge capacity while nearby hospitals had unused capacity. Not only is this inefficient, but it also could lead to a decreased quality of care at over-capacity hospitals.

**Methodology:** We develop and analyze a series of linear and mixed-integer programming models that solve variants of the demand and resource redistribution problem. We consider demand uncertainty and use robust optimization to ensure solution feasibility. We also incorporate a range of operational constraints and costs that decision makers may need to consider when implementing such a scheme.

**Results:** Our models are validated retrospectively using COVID-19 hospitalization data from New Jersey, Texas, and Miami, yielding at least an 85% reduction in required surge capacity relative to the observed outcome in each case. We show that such solutions can be operationally feasible and sufficiently robust against demand uncertainty.

**Managerial implications:** This work provides decision makers in healthcare systems with a practical

and flexible tool to reduce the surge capacity necessary to properly care for patients in cases when some facilities are over capacity.

**Keywords** COVID-19 Pandemic · Resource Allocation · Hospital Operations · Patient Transfers · Linear Programming · Mixed-Integer Optimization · Robust Optimization

## 1 Introduction

Since the first confirmed case of a SARS-CoV-2 infection in Wuhan, China, in January 2020, COVID-19 has rapidly evolved into a global pandemic (Hui et al. 2020). As of September 2020, there have been more than 32 million reported cases worldwide with the United States reporting the most confirmed COVID-19 cases among all countries at nearly 7 million. The world has seen more than 981 thousand deaths (Dong et al. 2020) as the result of the pandemic. COVID-19 has created an enormous burden on the capacities of healthcare systems across the globe, especially in Intensive Care Units (ICUs) (IHME et al. 2020). These healthcare capacity concerns have prompted state and local governments to intervene by instituting widespread closures and stay-at-home orders (Mervosh et al. 2020). Yet as the number of cases continues to soar around the globe and in the United States (US), healthcare capacity remains a concern at the hot spots of the pandemic.

In this paper, we present a framework to match the demand and available resources in healthcare networks. We, in particular, focus on scenarios in which individual hospitals in the network are under capacity stress during the current COVID-19 pandemic. In this framework, we develop several allocation optimization models to achieve a system-level load-balance by redistribution of patients (as opposed to resources such as ventilators or nurses) among different hospitals. The optimization models aim to minimize patient overflow in each hospital while also considering the operational constraints of patients transfers. We tested our proposed models on three case studies in New Jersey, Florida, and Texas using publicly available data gathered from state and healthcare agencies. We are currently in the process of adapting our models and tools for implementation in a large academic healthcare system in the US as it prepares for a potential new wave of SARS-CoV-2 infection in fall 2020. Our models will be used for optimal patient transfers within this healthcare system as well as informing strategic decisions on COVID-19 capacity management at each hospital.

*Motivation and Impact.* With the arrival of the pandemic to the US and the strain on healthcare capacity, the healthcare system responded by both reducing demand through canceling non-urgent procedures and also increasing capacity through creating new COVID-19-suitable beds as part of the guidelines from the Centers for Disease Control and Prevention (CDC) (CDC 2020a). While these measures have proven effective, delays in non-urgent care may result in adverse effects for patients. Increasing COVID-19 capacity is costly, slow, and may not be feasible at all healthcare centers. At the same time, healthcare centers have made efforts to optimize the use of complementary resources, including extending the use of Personal Protective Equipment (PPE) (CDC 2020b). However, in practice, these methods are often insufficient or unsustainable. Another compounding organizational drawback to conventional healthcare center responses is that often they are organized at individual healthcare centers level rather than system levels.

On a local level, patients tend to choose healthcare centers based on reputation or distance, leading to unbalanced patient loads across healthcare centers, which in turn decreases the overall quality of patient care (Varkevisser et al. 2012, Dreys 2013). A similar phenomenon can be observed on a larger scale where different regions experience varying pandemic-related demands at different times. Hence, while some healthcare centers may be operating at capacity, others might have additional capacity to spare. Current forecasts show that COVID-19 burdens healthcare centers far beyond their current capacity, especially in ICUs, which calls for a model that optimally distributes patients across healthcare centers to minimize overall capacity and resource shortages. Addressing the issue of unbalanced patient loads requires considering healthcare centers at the system level – across hospitals, counties, and even states – where there is potential for capacity pooling to lead to more efficient resource use (Vanberkel et al. 2012, Sims et al. 2015, Chod and Rudi 2005). Such approaches have been employed in ad-hoc manners during the pandemic. For instance, (Boudourakis

et al. 2020) describe patient transfer among a system of hospitals in New York despite operational challenges. Given the ongoing capacity concerns, healthcare systems are increasingly considering system-level interventions and patient transfers to maximize the utilization of available resources.

*Models and Case Studies.* In this work, we propose several optimization models to balance system loads in the face of an extreme stressor. Our models produce optimal and operationally feasible redistributions of newly admitted patients across healthcare systems to minimize total patient overflow and consequently reduce the strain on healthcare workers and infrastructure. The transfer of the patients happens at their point-of-entry to the healthcare center, which is often the Emergency Department (ED). This approach is in contrast with most current transfer approaches where a patient is transferred after admission to the first hospital, which may lead to increased inefficiency and length-of-stay (LOS). Although our reallocations are assumed to take place immediately upon patient arrival at a hospital, depending on EMS operational considerations, these transfers could potentially take place before arrival, rerouting patients on the way to hospitals.

The proposed models can not only redistribute patients (demand) but can also redistribute nurses and other critical resources across healthcare systems simultaneously, hence, further reducing shortage that may result from closed beds. We additionally consider various types of patients who require a different care path, e.g., acute-level care, ICU-level care, or any combination of them. This distinction plays an important operational role in capacity management as the available capacity, the patient demographic, and the LOS for each level of care is vastly different. These differences may become even more pronounced in subsequent waves of the pandemic. We also consider the uncertainties that exist in capacity management, especially during a novel pandemic, from the number of patients to the level of required care to their LOS. To address that, we consider a range of forecasts and distributions and propose a robust optimization model. It is worth noting that our models directly use the output of external patient count forecasts, which allows us flexibility in forecasting methodology. CDC (2020c) includes two instances of a situation where the states of Michigan and Texas opted to utilize such methods in extreme case surge intervals and were able to partially alleviate the situation.

Our methodology provides key insights from a managerial perspective. In a general setting, the recommendations from our model give centers experiencing excess demand a data-driven plan to reduce their loads and maintain adequate service quality. The flexibility afforded by these load reduction will also compound, allowing decision-makers to respond better to changing system demands. When combined with accurate forecasts, insights from our models can indicate vulnerable centers where an advanced expansion of capacity would be beneficial to the system as a whole.

Our results show that it was possible to reduce patient overflow by nearly 90% in New Jersey, one of the worst impacted states during the first wave of the pandemic. A challenge in our case studies was access to publicly and trustworthy available data. Therefore, we created a publicly available repository of relevant healthcare and forecasting data that can be used by other researchers and practitioners<sup>1</sup>.

A summary of the contributions of this work is as follows:

1. We provide a demand allocation optimization model to distribute patients across health care entities levels in the presence of extreme demand. This model allows for different patient care paths, each with an associated length of stay distribution for each bed type.
2. We show how to incorporate additional penalties and constraints to encourage the operational feasibility of solutions in the proposed model. We also propose a model that is robust against the most prominent sources of uncertainty in this setting.
3. We compile healthcare data relevant to our modeling and make it publicly available along with our model implementations and maintain an interactive website that allows users to specify model parameters and generate custom reallocation solutions.

---

<sup>1</sup><https://jhu-covid-optimization.github.io/covid-data/>

4. We present several case studies and our models are under implementation at a large academic health system in the US.

The rest of this paper is organized as follows. Related literature and the recent developments during the COVID-19 pandemic are presented in Section 2. Section 3 provides details on the proposed allocation models to improve load-balance in healthcare networks. Section 4 discusses the details of the data utilized in this study. The results of our analysis for the models and two case studies are illustrated in Section 5. Additional results (including a third case study), further details on the data preparation, and details on our robust model are provided in the accompanying online supplement. We conclude by discussing future directions in Section 6.

## 2 Related Work

In this section, we highlight some of the existing literature of different topics related to this work. We start by introducing some of the more prominent works in the general problem of resource allocation. Related literature on healthcare operations, modelings for pandemics and pandemic response are also provided. Finally, we discuss the existing works in the specific domain of healthcare resource allocation and the specific works related to the COVID pandemic. The general problem of resource allocation has been well studied in the literature, leading to the development of many different optimization methods in diverse domains. Details of such models can be found in the works of Hegazy (1999), Ren et al. (2018), Kuchuk et al. (2015), Cruz et al. (2004) and Tychogiorgos and Leung (2014). In particular, robust optimization approaches and applications are explored by Bertsimas and Sim (2004), Gabrel et al. (2014) and Najafi et al. (2013).

There is an extensive body of research around methods to optimize various healthcare system operations during periods that the system is not experiencing external stressors like shortages or pandemics. Among such works, (Luscombe and Kozan 2016) consider the problem of real-time dynamic emergency department (ED) scheduling with the goal of minimizing patient wait times. Similarly, Otegbeye et al. (2015) develop a tool which utilizes retrospective ED data to generate optimal nurse shift schedules. Litvak et al. (2008) present mathematical support for a system-wide pool of reserved emergency ICU beds to improve patient care and successful admission rates efficiently through network-wide cooperation. There have also been studies that consider high demand periods in their approaches too. For instance, Toner and Waldhorn (2006) emphasize that during high demand periods such as a pandemic, cooperation between different healthcare centers becomes crucial to accommodate extreme system-wide stress. Nonetheless, much of the existing literature in the optimization healthcare operations only considers the setting of a single healthcare center or does not consider high-demand settings where some healthcare centers rapidly hit their load capacity. For extensive surveys regarding operations research applications in healthcare systems, the reader is referred to the works of Papageorgiou (1978) and Rais and Viana (2011).

In the specific setting of the existence of a pandemic, considerable work has been done regarding optimal healthcare responses. Brandeau (2005) explore the problem of general resource allocation in response to a pandemic. Toner and Waldhorn (2006) develop a set of best practices for influenza pandemic preparedness and response, and Halpern and Miller (2020) Cooper et al. (2006) explore the challenges and limitations of increasing critical care supply and shutting down air travel respectively in the presence of a pandemic. Other works in the literature consider more granular levels of analysis. For instance, CDC (2020b) propose operational strategies to optimize the use of face-masks and other personal protective equipment and Mehrotra et al. (2020a) build a framework for rapidly moving patient primary care to a virtual modality. Some more recent studies have explored ethical approaches to rationing critical resource allocations during the COVID-19 pandemic (Emanuel et al. 2020, White and Lo 2020). Finally, Judson et al. (2020) create a patient self-scheduling and self-triage tool designed to decrease patient wait times during the COVID-19 pandemic. It should be noted that the above studies do not address the issue of unbalanced loads among different entities of a healthcare system and potential resource allocation approaches to alleviate such unbalanced demand on a healthcare systems under extreme stress, which is the primary attention of this work.

Attention to developing tools and methods to track and model pandemics have been central in the literature (Pei et al. 2018), and the literature has grown considerably since the emergence of the COVID-19 pandemic (Petroopoulos and Makridakis 2020, Jewell et al. 2020, Roda et al. 2020, Perc et al. 2020). Recently, Lewnard et al. (2020) and Rees et al. (2020) retrospectively analyze and synthesize COVID-19 healthcare data to provide key qualitative and quantitative insights into the pandemic’s nature and spread. Meanwhile, Dong et al. (2020) provide valuable data by tracking cases across all nations and in particular, in the United States at the county level, and Weissman et al. (2020) and IHME et al. (2020) provide forecasts for patient loads at local hospital level and national state level respectively. In general, existence of such data and tools is paramount in developing meaningful models for optimizing different operations during pandemics.

We close the related works section by pointing out existing works in the problem of resource allocation using central stockpiles. This problem has been previously studied by Arora et al. (2010) and Mehrotra et al. (2020b) in deterministic and stochastic settings respectively. More specifically, the issue of optimal vaccine allocation has been addressed by Longini Jr et al. (1978) and Lampariello and Sagratella (2020) in non-COVID and COVID settings respectively and operational considerations of COVID-19 patient transfers are discussed in CDC (2020c) and Liew et al. (2020). As a real world example of patient allocation schemes during the COVID-19 pandemic, Boudourakis et al. (2020) present the successes and challenges associated with patient transfers among New York City hospitals. In the most relevant existing literature, Bai and Zhang (2014), Sun et al. (2014) and Lacasa et al. (2020) discuss the problem of optimal patient allocation in a pandemic. Lacasa et al. (2020) consider the problem of distributing a single resource or demand across a regular geometric graph with healthcare center as vertices. They provide solutions as a set of resource transfers using random search optimization. However, it should be noted that such solutions are not guaranteed to be optimal, only correspond to a single time step, do not support complementary resources or secondary operational constraints, and are not robust against any particular data uncertainties. On the other hand, Bai and Zhang (2014) cast the problem of patient load-balancing in a pandemic setting as a max-flow problem, offering a solution containing a set of patient incentives to optimally balance healthcare center loads. However, rather than considering minimizing the capacity overflow, they seek to minimize patient cost, suggesting that their model may be more applicable in a lower-demand setting. Finally, Sun et al. (2014) produce a model similar to the models in this work, however, similar to Bai and Zhang (2014), they attempt to minimize patient travel distance and assume a constant length of stay for each group of patients. They also provide fewer secondary operational constraints to improve practical solution feasibility. In this work, we extend the capabilities of the models in the literature by proposing a new modeling methodology for the optimum allocation of patients and/or complementary resources subject to a variety of secondary operational constraints and potential uncertainties in different parameters like the number of admitted patients, beds and supplementary resources.

### 3 Methodology

In this section, we present a series of linear optimization (LP and MILP) models to solve the load-balancing problem in the presence of an extreme stressor to the system. The healthcare system is modeled as a group of nodes, each with a certain capacity accommodating COVID-19 patients. These nodes are put under stress at different points in time by the surges in demand happening in the course of the pandemic. To address this problem, we formulate a linear optimization problem that leverages the whole system’s capacity to help individual nodes at surge times by transferring patients between the nodes. Our secondary goal is to adapt the models into probable existing constraints in the problem domain to make them more flexible and feasible to implement. Finally, the models are modified to make the resulting capacity-sharing scheme robust against uncertainties in the predictions or the parameters.

We start with a simplified scenario where the nodes’ capacities are governed by a single resource type, staffed beds for COVID-19 patients. The introduced formulation is further developed in Section 3.1.2 to allow for different patient care path groups with distinct length of stay (LOS) distributions, care requirements and bed types, e.g., ICU and non-ICU.

In the next step, to distribute patients more realistically and to address operational concerns, various optional penalties and constraints are introduced, which can be added to the optimization scheme as needed.

Next, we propose a complementary model in which, alongside the primary load-balancing (transferring patients in our study), certain resource types can also be redistributed among the nodes. In this work, nurses are considered as an example of transferable resources. We choose nurses as primary transferable resources since they are as crucial as the number of beds in governing the nodes' capacities and are easily transferable (Section 3.2). Finally, we propose a robust optimization model to account for the uncertainties in the daily numbers of admitted COVID-19 patients and available beds.

In this paper, the common terminology is specific to healthcare system load-balancing; however, it is noteworthy that the methodologies outlined here are not limited to healthcare domains during the COVID-19 pandemic. Any service industry with challenges regarding potential stressors and resource allocation can calibrate the models to their specific needs.

### 3.1 Demand Allocation Models

In this section, we develop an optimization scheme to re-distribute demand between nodes so that the total surge capacity required to accommodate the demand is minimized. We specifically consider the case where demand is COVID-19 patients and capacity is beds, but the model is flexible enough to be used for other applications as well. In this case, each node is assumed to have a certain fixed number of staffed beds that meet appropriate clinical requirements to care for COVID-19 patients, which compose the node's capacity. COVID-19 patients are admitted to the nodes on a daily basis, composing the demand for staffed beds. The capacity and number of admissions are assumed to be known, as is discussed in Section 4 and the online supplement for this paper. We utilize linear models to first develop our base formulation for patient allocation in Section 3.1.1. This model is then extended to enable it to distinguish between various groups of patients who have a different care path (e.g., patients who go to the ICU at some point in their care and patients who do not) in Section 3.1.2. We then introduce several additional constraints and penalties which increase the model's practical utility and flexibility in Section 3.1.3. The details of our assumptions, parameters, variables and other notation for all models are summarized in Table 1 at the end of this section.

#### 3.1.1 Patient Allocation Formulation.

We begin with formulating a linear optimization model to transfer COVID-19 patients between nodes in a given healthcare system consisting of nodes with fixed capacities governed by a single resource type, namely staffed beds, with the objective to minimize the total patient overflow across the system over the considered period, subject to basic feasibility constraints, as follows:

$$\underset{\omega}{\text{minimize}} \quad \sum_{i \in \mathcal{N}} \sum_{t \in \mathcal{T}} \omega_{i,t} \quad (1a)$$

$$\text{subject to} \quad \sum_{j \in \mathcal{N}} s_{i,j,t} \leq p_{i,t} \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (1b)$$

$$\alpha_{i,t} - b_i \leq \omega_{i,t} \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (1c)$$

$$\begin{aligned} \alpha_{i,t} = & (p_{i,0} - \sum_{t'=1}^t d_{i,t'}) + \sum_{j \in \mathcal{N}} s_{i,j,t} \\ & + \sum_{t'=1}^t \{[1 - \mathcal{L}(t - t')][p_{i,t'} + \sum_{j=1}^N (s_{j,i,t'} - s_{i,j,t'})]\} \end{aligned} \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (1d)$$

$$s_{i,j,t} = 0 \quad \forall (i,j) \in \overline{E(G)}, t \in \mathcal{T} \quad (1e)$$

$$s_{i,j,t} \geq 0, \quad \omega_{i,t} \geq 0 \quad \forall i,j \in \mathcal{N}, t \in \mathcal{T} \quad (1f)$$

In this formulation, later referred to as the base model, we assume that only newly admitted patients may be transferred between nodes, which is enforced by constraint (1b), ensuring that the number of patients transferred away does not exceed the number of newly admitted patients  $p_{i,t}$ . The overflow for node  $i$  at time  $t$ , denoted by  $\omega_{i,t}$ , is defined by constraint (1c) in conjunction with the constraint that  $\omega_{i,t} \geq 0$ . It is also assumed that transfers between certain nodes can be infeasible (e.g., based on distance or existing relationships between hospitals). These relationships and restrictions can be formally modeled as a graph  $G$  where all healthcare centers are considered as vertices, and there exists an edge from node  $i$  to node  $j$  if patients can be transferred between them. This assumption is asserted by constraint (1e). Expression (1d) represents the number of active patients at node  $i$  at time  $t$ . The first term of this expression captures the number of remaining initial patients. The third term incorporates the cumulative patient length of stay distribution  $\mathcal{L}$  to capture the number of admitted and reallocated patients remaining at node  $i$  at time  $t$ . Since patient transfer between care centers typically requires resources from both centers on the day of the transfer, we count the transferred patients as active at both nodes  $i$  and  $j$  at time  $t$ , which is captured by the second term of the expression. More details about the parameters and variables of this formulation are provided in Table 1.

This model assumes a single care path for COVID-19 patients, and a single limiting capacity, total staffed beds. In the next section, we expand the base model to take into account patients with different care paths and their varying resource requirements.

### 3.1.2 Care Path Group Patient Allocation Formulation.

In this section, we extend the base LP model discussed in Section 3.1.1 to provide a model capable of considering a setting where instead of a single type of bed and patients, a set of patient groups  $\mathcal{G}$  and also multiple bed types  $\mathcal{B}$  are present. Each patient is assigned to a specific patient group at admission time. During the course of the treatment of a patient, they may undergo different care paths, meaning that they can be transferred across the patient groups and potentially need various bed types during the different phases of their treatments.

To model the care paths, directed graph  $G_{group}$  is defined such that each node of it corresponds to a specific phase of treatment with its own resource type requirements, which is referred to as a patient group, and the edges connecting the consequent patient groups. In this setting,  $G_{group}$  is assumed to be a disjoint union of directed trees with edges directed toward their own roots (also known as an in-forest graph). This assumption is asserted to disallow cycles in a patient group transfer scheme and simultaneously require each patient to have a unique transfer path until they are discharged.

For formulation purposes, a function  $f$  mapping  $\mathcal{G}$  to  $\mathcal{B}$  is defined, implying that each patient group is associated with exactly one bed type.

$$\underset{\omega}{\text{minimize}} \quad \sum_{\beta \in \mathcal{B}} \sum_{i \in \mathcal{N}} \sum_{t \in \mathcal{T}} \omega_{\beta,i,t} \quad (2a)$$

$$\text{subject to} \quad \sum_{j \in \mathcal{N}} s_{g,i,j,t} \leq p_{g,i,t} \quad \forall g \in \mathcal{G}, i \in \mathcal{N}, t \in \mathcal{T} \quad (2b)$$

$$\sum_{g \in \text{img}\{f^{-1}(\beta)\}} (\alpha_{g,i,t} + \sum_{j \in \mathcal{N}} s_{g,i,j,t}) - b_{\beta,i} \leq \omega_{\beta,i,t} \quad \forall \beta \in \mathcal{B}, i \in \mathcal{N}, t \in \mathcal{T} \quad (2c)$$

$$\alpha_{g,i,t} = p_{g,i,0} + \sum_{t'=1}^t (\chi_{g,i,t'} - \gamma_{g,i,t'}) \quad \forall g \in \mathcal{G}, i \in \mathcal{N}, t \in \mathcal{T} \quad (2d)$$

$$\chi_{g,i,t} = p_{g,i,t} + \sum_{g': g' \sim g} \gamma_{g',i,t} + \sum_{j \in \mathcal{N}} (s_{g,j,i,t} - s_{g,i,j,t}) \quad \forall g \in \mathcal{G}, i \in \mathcal{N}, t \in \mathcal{T} \quad (2e)$$

$$\gamma_{g,i,t} = d_{g,i,t} + \sum_{t'=1}^t \ell_g(t-t') \chi_{g,i,t'} \quad \forall g \in \mathcal{G}, i \in \mathcal{N}, t \in \mathcal{T} \quad (2f)$$

$$s_{g,i,j,t} = 0 \quad \forall g \in \mathcal{G}, (i,j) \in \overline{E(G)}, t \in \mathcal{T} \quad (2g)$$

$$s_{g,i,j,t}, \omega_{\beta,i,t} \geq 0 \quad \forall g \in \mathcal{G}, \beta \in \mathcal{B}, i, j \in \mathcal{N}, t \in \mathcal{T} \quad (2h)$$

In the above model, constraints (2b) through (2h) are generalizations of the constraints (1b) through (1f) from the base model introduced in Section (3.1.1). Expression (2d) represents the number of active patients in group  $g$  at node  $i$  at time  $t$ . The first term captures initial patients while the second term accounts for the sum of net active patient changes in group  $g$  at node  $i$ . Expression (2e) represents the number of patients entering group  $g$  in node  $i$  at time  $t$ . The first term captures admitted patients into group  $g$ , the second term sums patients leaving other groups that transfer to the group  $g$ , and the third term includes net patient transfers in group  $g$  at node  $i$  at time  $t$ . Finally, expression (2f) represents the number of patients leaving group  $g$  in node  $i$  at time  $t$ . The first term of Expression (2f) captures initial patients discharged while the second term calculates the number of patients leaving group  $g$  for another group, at node  $i$  at time  $t$ . Notice that although expressions (2d) through (2f) immediately constitute a closed form expression for the number of active patients, they provide a method for computing one recursively which is guaranteed to terminate by our requirement that  $G_{\text{group}}$  be a in-forest. More details on the parameters and variables of this formulation are provided in Table 1.

The care path group patient allocation model developed here is capable of minimizing patient overflow across multiple patient groups and bed types in a given system. In the following section we provide a collection of additional constraints and penalties that can be added to the so far developed scheme, which add flexibility and capability to the models to account for probable operational considerations in the optimal patient allocation problem.

### 3.1.3 Optional Constraints and Considerations.

We now build on the group patient LP formulation to add additional penalties and constraints to ensure the operational feasibility of solutions. In particular, we have added terms to penalize undesirable solutions of the model presented in



3.1.2. Details on the additional variables and parameters are provided in Table 1.

$$\begin{aligned} \text{penalty} = & C_{\text{sent}} \sum_{g \in \mathcal{G}} \sum_{i,j \in \mathcal{N}} \sum_{t \in \mathcal{T}} s_{g,i,j,t} + C_{\text{smooth}} \sum_{i,j \in \mathcal{N}} \sum_{t \in \mathcal{T}} \delta_{i,j,t} \\ & + C_{\text{balance}} \sum_{\beta \in \mathcal{B}} \sum_{i \in \mathcal{N}} \sum_{t \in \mathcal{T}} \phi_{\beta,i,t} + C_{\text{setup}} \sum_{i,j \in \mathcal{N}} \rho_{i,j} \end{aligned} \quad (3a)$$

$$-\delta_{i,j,t} \leq \sum_{g \in \mathcal{G}} (s_{g,i,j,t-1} - s_{g,i,j,t}) \leq \delta_{i,j,t} \quad \forall i, j \in \mathcal{N}, t \in \mathcal{T} \setminus \{1\} \quad (3b)$$

$$\frac{1}{b_{\beta,i}} \sum_{g \in \text{img}\{f^{-1}(\beta)\}} \alpha_{g,i,t} - R_{\text{thresh}} \leq \phi_{\beta,i,t} \quad \forall \beta \in \mathcal{B}, i \in \mathcal{N}, t \in \mathcal{T} \quad (3c)$$

$$M \sum_{g \in \mathcal{G}} \sum_{t \in \mathcal{T}} (s_{g,i,j,t} + s_{g,j,i,t}) \geq \rho_{i,j} \quad \forall i, j \in \mathcal{N}, j > i \quad (3d)$$

$$m \sum_{g \in \mathcal{G}} \sum_{t \in \mathcal{T}} (s_{g,i,j,t} + s_{g,j,i,t}) \leq \rho_{i,j} \quad \forall i, j \in \mathcal{N}, j > i \quad (3e)$$

$$M \sum_{g \in \mathcal{G}} \sum_{j \in \mathcal{N}} s_{g,i,j,t} \geq \nu_{k,i,t} \quad \forall k \in \{1, 2\}, i \in \mathcal{N}, t \in \mathcal{T} \quad (3f)$$

$$m \sum_{g \in \mathcal{G}} \sum_{j \in \mathcal{N}} s_{g,i,j,t} \leq \nu_{k,i,t} \quad \forall k \in \{1, 2\}, i \in \mathcal{N}, t \in \mathcal{T} \quad (3g)$$

$$\nu_{k,i,t} + m \sum_{g \in \mathcal{G}} \sum_{j \in \mathcal{N}} \sum_{t'=t}^{\min\{t+T_{\text{switch}}, T\}} s_{g,j,i,t'} \leq 1 \quad \forall k \in \{1, 2\}, i \in \mathcal{N}, t \in \mathcal{T} \quad (3h)$$

$$\sum_{g \in \text{img}\{f^{-1}(\beta)\}} \alpha_{g,i,t} \leq \max b_{\beta,i} \quad \forall \beta \in \mathcal{B}, i \in \mathcal{N}, t \in \mathcal{T}, \alpha'_{\beta,i,t} \leq b_{\beta,i} \quad (3i)$$

$$\sum_{g \in \text{img}\{f^{-1}(\beta)\}} \alpha_{g,i,t} \leq \alpha'_{\beta,i,t} \quad \forall \beta \in \mathcal{B}, i \in \mathcal{N}, t \in \mathcal{T}, \alpha'_{\beta,i,t} > b_{\beta,i} \quad (3j)$$

$$\sum_{g \in \mathcal{G}} s_{g,i,j,t} \in \{0\} \cup [S_{\min}, \infty), \quad \delta_{i,j,t}, \phi_{\beta,i,t} \geq 0 \quad \forall \beta \in \mathcal{B}, i, j \in \mathcal{N}, t \in \mathcal{T} \quad (3k)$$

$$\rho_{i,j}, \nu_{1,i,t}, \nu_{2,i,t} \in \{0, 1\} \quad \forall i, j \in \mathcal{N}, t \in \mathcal{T} \quad (3l)$$

$$\alpha'_{\beta,i,t} = \sum_{g \in \text{img}\{f^{-1}(\beta)\}} [p_{g,i,0} + \sum_{t'=1}^t (\chi'_{g,i,t'} - \gamma'_{g,i,t'})] \quad \forall \beta \in \mathcal{B}, i \in \mathcal{N}, t \in \mathcal{T} \quad (3m)$$

$$\chi'_{g,i,t} = p_{g,i,t} + \sum_{g': g' \sim g} \gamma'_{g',i,t} \quad \forall g \in \mathcal{G}, i \in \mathcal{N}, t \in \mathcal{T} \quad (3n)$$

$$\gamma'_{g,i,t} = d_{g,i,t} + \sum_{t'=1}^t \ell_g(t-t') \chi'_{g,i,t'} \quad \forall g \in \mathcal{G}, i \in \mathcal{N}, t \in \mathcal{T} \quad (3o)$$

In the above formulation, the penalty represents a sum of terms that can be added to the group patient allocation formulation from section (3.1.2) along with their corresponding constraints. The first cost coefficient,  $C_{\text{sent}}$ , penalizes the objective function for the total number of patients sent by the solution.  $C_{\text{smooth}}$  penalizes for the absolute value of the difference between patient transfer quantities from node  $i$  to node  $j$  at time  $t-1$  and time  $t$ . The corresponding variable  $\delta_{i,j,t}$  is determined by constraints (3b). These constraints are motivated by the idea that a more consistent transfer rate from node  $i$  to node  $j$  may be more operationally feasible.  $C_{\text{balance}}$  encourages more balanced patient loads by penalizing the objective function when patient load at a node exceeds the balancing threshold ratio,  $R_{\text{thresh}}$ . Constraint (3c) defines the corresponding dummy variable  $\phi_{\beta,i,t}$ . Finally, the last term in the penalty function including  $C_{\text{setup}}$  for the first time a transfer takes place between nodes  $i$  and  $j$ , accounting for overhead cost of establishing a transfer relationship. Constraints (3d) and (3e) define the  $\rho_{i,j}$  binary variable accordingly.

Constraints (3f) through (3h) enforce a minimum gap of  $T_{\text{switch}}$  days between when a node may send and receive patients. This is done by introducing the binary variables  $\nu_{1,i,t}$  and  $\nu_{2,i,t}$  and using them to constrain the sent and received patients for the next  $T_{\text{switch}}$  days. These constraints contribute to smoothness as well, ensuring nodes will be unable to switch quickly between sending patients and receiving them (or do both at once). Then, Constraint (3i) ensures that no node experiences a patient overflow at a time when they otherwise would not, and Constraint (3j) ensures that no node overflow is made more severe by our reallocation. Next, Constraint (3k) contains a non-convex constraint that can be implemented using binary variables which guarantees that if a patient transfer occurs, at least  $S_{\min}$  patients are sent.

Expressions (3m) through (3o) are analogous to Expressions (2d) through (2f), representing active patients, entering patients, and exiting patients for each group, except that these new expressions do not consider patient reallocation and  $\alpha'_{\beta,i,t}$  is aggregated by bed type. Therefore Expression (3m) gives the number of active patients in bed type  $\beta$  at node  $i$  at time  $t$  without any patient reallocation. So far we have developed a patient transferal scheme with considerations for variety in care paths and bed types. In the next section this scheme is further expanded to do a simultaneous resource redistribution along with the patient transfers.

### 3.2 Combined Demand and Resource Allocation Model

In this section we build on the group patient LP formulation from section 3.1.2 to allocate nurses as a primary transferable resource, along with the patients. In order for a patient to receive proper treatment we require that they both have a bed and adequate nurse care. First, a nurse allocation LP model is introduced. Secondly, a set of nurse specific optional constraints are added, addressing the issue of artificial shortage of static supply for nurses (compared to patients' dynamic demand). It is noteworthy that to simultaneously allocate patients and nurses, we can combine the variables and constraints from the models in this section with the variables and constraints from Section 3.1 and take our new objective function to be a weighted sum of all of the individual model objective functions.

#### 3.2.1 Nurse Allocation Formulation.

Nurses are a primary resource for caring for COVID-19 patients and so can certainly affect the nodes' capacities. They also have the capability to transfer across the nodes as needed. The following formulation is used to balance the number of nurses among the nodes.

$$\underset{\theta}{\text{minimize}} \quad \sum_{i \in \mathcal{N}} \sum_{t \in \mathcal{T}} \theta_{i,t} \quad (4a)$$

$$\text{subject to} \quad \sum_{j \in \mathcal{N}} \sigma_{i,j,t} \leq \eta_{i,t} \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (4b)$$

$$q_{i,t} - \eta_{i,t} \leq \theta_{i,t} \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (4c)$$

$$\eta_{i,t} = n_i + \sum_{j=1}^N \sum_{t'=ts(t)}^t (\sigma_{j,i,t'} - \sigma_{i,j,t'}) \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (4d)$$

$$q_{i,t} = \sum_{\beta \in \mathcal{B}} \sum_{g \in \text{img}\{f^{-1}(\beta)\}} Q_{\beta} \alpha_{g,i,t} \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (4e)$$

$$\sigma_{i,j,t} = 0 \quad \forall (i,j) \in \overline{E(G)}, t \in \mathcal{T} \quad (4f)$$

$$\sigma_{i,j,t} \geq 0, \theta_{i,j} \geq 0 \quad \forall i,j \in \mathcal{N}, t \in \mathcal{T} \quad (4g)$$

Equations (4b) through (4g) are analogous to equations (1b) through (1f) in section (3.1). Expression (4d) represents the number of active nurses and expression (4e) represents the nurse demand which is taken to be proportional to the number of active patients in group  $g$  at node  $i$  at time  $t$ . When allocating patients and nurses simultaneously,  $\alpha_{g,i,t}$  will

depend on the model's patient allocation. In the following section we provide additional constraints which can be added to the nurse allocation model to produce more operationally feasible solutions.

### 3.2.2 Optional Constraints and Considerations.

In this section we provide a set of constraints designed to account for the fact that if a node is experiencing a shortage of nurses, either after or during patient allocation, then that node should not have any nurses deployed to other nodes.

$$n_i \leq \eta_{i,t} \quad \forall i \in \mathcal{N}, t \in T, q_{i,t} \geq n_{i,t} \quad (5a)$$

$$m(q_{i,t} - \eta_{i,t}) \leq \kappa_{i,t} \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (5b)$$

$$1 + m(q_{i,t} - \eta_{i,t}) \geq \kappa_{i,t} \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (5c)$$

$$\eta_{i,t} \geq n_i \quad \forall i \in \mathcal{N}, t \in T, \kappa_{i,t} = 1 \quad (5d)$$

$$\kappa_{i,t} \in \{0, 1\} \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (5e)$$

Equations (5b), (5d) and (5e) enforce that for any day that a region has a shortage of nurses, it has at least its initial supply of nurses. Optional penalties and constraints from 3.1.2 can also be applied to nurses, but we forgo rewriting them here for simplicity.

### 3.3 Robust Optimization Model

Thus far our models have all implicitly relied on the assumption that we know the input data with certainty. In practice, however, regardless of whether our methodology is implemented prospectively or retrospectively, there is almost surely some extent of uncertainty in the data. We therefore construct a Robust Optimization model which ensures the optimal solution remains feasible for all scenarios in the uncertainty set we consider. The model we present here addresses uncertainty on the number of admitted patients. The number of admitted patients is particularly impactful on the solution because it governs the number of active patients that need to be accommodated and it limits the number of patients that can be transferred. We also expect that this input will be the most uncertain data for prospective studies since we must rely on forecasting to generate the predictions.

We start with a box uncertainty set, meaning that we consider an upper and lower bound on the number of admitted patients, and assume that the true number of admitted patients must fall somewhere within this range. Specifically,  $p_{i,t} \in P_{i,t} = [\bar{p}_{i,t} - \tilde{p}_{i,t}^-, \bar{p}_{i,t} + \tilde{p}_{i,t}^+] \forall i \in \mathcal{N}, t \in \mathcal{T}$ . We define  $\bar{p}_{i,t}$  to be the nominal value of the number of admitted patients, which is what is used in the preceding models. Additionally, we assume that  $\tilde{p}_{i,t}^-, \tilde{p}_{i,t}^+ \geq 0$  and  $\tilde{p}_{i,t}^- \leq \bar{p}_{i,t}$ . Note that a distribution over the number of admitted patients is not specified. Rather, this method robustifies against any distribution with a fixed upper and lower bound, making it highly flexible.

Previous work in Robust Optimization, however, including Bertsimas and Sim (2004) have shown that such models can be overly conservative with a high "price of robustness" – that is they suffer a large increase to the optimal objective function value as compared with the certain optimization model. We therefore adopt an "uncertainty budget" inspired by Bertsimas and Sim (2004). Instead of assuming that the uncertain values can be anywhere between the upper and lower bound we constrain the uncertain values to be equal to the nominal value except on a parameter  $\Gamma$  number of days, when it deviates to the upper or lower bound. Formally,  $p_{i,t} = \bar{p}_{i,t} + \min\{0, \xi_{i,t} \tilde{p}_{i,t}^-\} + \max\{0, \xi_{i,t} \tilde{p}_{i,t}^+\}$  such that  $\xi_i \in \{0, 1\}^t, \|\vec{\xi}_i\|_1 = \Gamma \forall i \in \mathcal{N}$ . This formulation captures the intuition that the true value of the number of admitted patients is generally equal to the nominal value, but will deviate on some days. The practitioner may vary  $0 \leq \Gamma \leq |\mathcal{N}|$  where at the extremes the fully robust and certain models are recovered respectively. Note that in practice we can relax the constraint  $\xi_i \in \{0, 1\}^t$  to  $\|\vec{\xi}_i\|_\infty \leq 1$  and our solution method will recover the original constraint. Note also that due to the uncertainty budget, the number of admitted patients each day for a given node are dependent on one another.

Consequently, a realization from the uncertainty set for a given node must consist of a sequence of admitted patients per day for that node.

A motivating factor for employing robust optimization and this uncertainty set in particular is that the robust model may be formulated as an LP with the same number of constraints and variables as the equivalent non-robust model, meaning that it can be solved efficiently in theory and practice. The details of this re-formulation may be found in the online supplement to this paper.

The following is the base patient allocation model presented in Section 3.1.1 made robust against our identified uncertainty set.

$$\underset{\omega}{\text{minimize}} \quad \sum_{i \in \mathcal{N}} \sum_{t \in \mathcal{T}} \omega_{i,t} \quad (6a)$$

$$\text{subject to} \quad \sum_{j \in \mathcal{N}} s_{i,j,t} \leq p_{i,t} \quad \forall i \in \mathcal{N}, t \in \mathcal{T}, p_{i,t} \in P_{i,t} \quad (6b)$$

$$\alpha_{i,t} - b_i \leq \omega_{i,t} \quad \forall i \in \mathcal{N}, t \in \mathcal{T}, \alpha_{i,t} \in A_{i,t} \quad (6c)$$

$$\begin{aligned} A_{i,t} = & \left\{ (p_{i,0} - \sum_{t'=1}^t d_{i,t'}) + \sum_{j \in \mathcal{N}} s_{i,j,t} \right. \\ & + \sum_{t'=1}^t ([1 - \mathcal{L}(t - t')][p_{i,t'} + \sum_{j=1}^N (s_{j,i,t'} - s_{i,j,t'})]) \\ & : p_{i,t'} = \bar{p}_{i,t'} + \min \{ 0, \xi_{i,t'} \bar{p}_{i,t'}^- \} + \max \{ 0, \xi_{i,t'} \bar{p}_{i,t'}^+ \} \\ & \left. \forall t' \in \{ 1, \dots, t \}, \|\xi_{i,1:t}\|_{\infty} \leq 1, \|\xi_{i,1:t}\|_1 \leq \Gamma \right\} \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (6d) \end{aligned}$$

$$s_{i,j,t} = 0 \quad \forall (i, j) \in \overline{E(G)}, t \in \mathcal{T} \quad (6e)$$

$$s_{i,j,t} \geq 0, \quad \omega_{i,t} \geq 0 \quad \forall i, j \in \mathcal{N}, t \in \mathcal{T} \quad (6f)$$

This initial robust formulation is potentially very large relative to the base patient allocation model and is computationally heavy. Assuming that the uncertainty set  $P$  is finite, constraints (6b) and (6c) each represent  $|\mathcal{N}| \cdot |\mathcal{T}| \cdot |P|$  constraints. In comparison, the base model has  $|\mathcal{N}| \cdot |\mathcal{T}|$  constraints, which makes the robust model possibly much more time consuming to solve. If  $P$  is instead infinite, then the robust model will have infinite number of constraints. However, it is possible to reformulate the robust model in such a way that it contains the same number of constraints and variables as the base model.

In all the models developed in this work, the number of admitted patients is a constant term in constraints rather than a coefficient on some variable. This means that the uncertainty only exists on the right-hand-side of the linear programming formulation when it is put in standard form, which makes the robust model simpler and easier to solve than a general robust optimization model. In this case, to ensure that a constraint remains feasible for all realizations in the uncertainty set, we can replace each uncertain constraint with a certain constraint using the worst-case from the uncertainty set of the parameter to compute the right-hand-side. If the constraint is of the form  $ax \leq b$  where  $b$  is in some uncertainty set  $B$  then the constraint is feasible for all  $b \in B$  if and only if it is feasible for  $b^- = \min B$ . Similarly, if the constraint is of the form  $ax \geq b$  where  $b$  is in some uncertainty set  $B$  then the constraint is feasible for all  $b \in B$  if and only if it is feasible for  $b^+ = \max B$ .

The worst case for constraint (6c) in the initial robust formulation is the realization from the uncertainty set that yields the largest number of active patients. Therefore, to solve the model, it must be determined which realization yields the largest number of active patients under the constraint  $\|\xi_i\| \leq \Gamma$ . For a given node  $i$  and day  $t$  the active patients is a value from the set  $A_{i,t}$  defined in constraint (6d). This expression is maximized, generating the worst case, when  $\sum_{t'=1}^t \xi_{i,t'}(1 - \mathcal{L}(t - t'))\bar{p}_{i,t'}^+$  is maximized since all other terms are constant or non-positive. Assuming that each

$(1 - \mathcal{L}(t - t'))\tilde{p}_{i,t'}^+$  is unique,  $\xi_{i,t}$  will be 1 for the  $\Gamma$  days with largest  $(1 - \mathcal{L}(t - t'))\tilde{p}_{i,t'}^+$  and zero otherwise. This results in a closed-form expression for  $\xi_{i,t}$  that produces the worst case right-hand-side for constraint (6c). Notably, this allows us to solve the robust model without adding any additional constraints or variables to the base LP model.

However, an important consideration for our robust model is that the worst-case scenario from the uncertainty set will not be consistent between constraints. We therefore must select a worst-case scenario per-constraint. This means that the robust model is more conservative than suggested by the uncertainty budget, but in practice  $\Gamma$  can still be selected such that the model has the desired level of robustness.

The reformulation of the robust model is therefore as follows:

$$\text{minimize}_{\omega} \quad \sum_{i \in \mathcal{N}} \sum_{t \in \mathcal{T}} \omega_{i,t} \quad (7a)$$

$$\text{subject to} \quad \sum_{j \in \mathcal{N}} s_{i,j,t} \leq p_{i,t}^-, \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (7b)$$

$$\alpha_{i,t}^+ - b_i \leq \omega_{i,t}, \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (7c)$$

$$s_{i,j,t} = 0 \quad \forall (i, j) \in \overline{E(G)}, t \in \mathcal{T} \quad (7d)$$

$$s_{i,j,t} \geq 0 \quad \forall i, j \in \mathcal{N}, t \in \mathcal{T} \quad (7e)$$

$$\omega_{i,t} \geq 0 \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (7f)$$

$$p_{i,t}^- = \bar{p}_{i,t} - \tilde{p}_{i,t}^- \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (7g)$$

$$p_{i,t}^+ = \bar{p}_{i,t} + \tilde{p}_{i,t}^+ \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (7h)$$

$$\begin{aligned} \alpha_{i,t}^+ &= (p_{i,0} - \sum_{t'=1}^t d_{i,t'}) \\ &+ \sum_{t'=1}^t \{ [1 - \mathcal{L}(t - t')] [(1 - \zeta_{i,t,t'})\bar{p}_{i,t'} + \zeta_{i,t,t'}\tilde{p}_{i,t'}^+] \\ &+ \sum_{j=1}^N (s_{j,i,t'} - s_{i,j,t'}) \} \end{aligned} \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (7i)$$

Where:

$$\zeta_{i,t} = \arg \max_{\hat{\zeta}_{i,t}} \sum_{t'=1}^t \mathcal{L}(t - t') \tilde{p}_{i,t'}^+ \hat{\zeta}_{i,t,t'} \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (8a)$$

$$\text{subject to} \quad \hat{\zeta}_{i,t} \in \{0, 1\}^t \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (8b)$$

$$\|\hat{\zeta}_{i,t}\|_1 = \min \{ \Gamma, t \} \quad \forall i \in \mathcal{N}, t \in \mathcal{T} \quad (8c)$$

### 3.4 Solution Evaluation

We consider the overflow in patient-days, which is defined as the number of patients minus the number of beds, summed over all node-day pairs such that the number of patients is greater than the number of beds:  $\sum_{i \in \mathcal{N}} \sum_{t \in \mathcal{T}} \max \{ 0, \alpha_{i,t} - b_i \}$ . The overflow is the amount of extra capacity that a node will have to create to properly care for all its patients, or else it must provide sub-standard care or turn patients away. Each of the models directly minimizes this key metric. A second metric, overflow reduction, is the decrease in total overflow under the model's transfer scheme as a percentage of the baseline overflow, which is the overflow assuming no transfers were made. Additionally, the size of individual node-day overflows is considered as large overflows on a given node-day necessitate a large amount of surge capacity to avoid inadequate patient care. This means good solutions will have smaller individual node-day overflows, even if this means the total overflow is spread out over more node-days. We

therefore evaluate models on the mean, median, and maximum non-zero overflow. Overflow can also be computed system-wide rather than for a given node-day pair by aggregating all patients and capacity. This system-wide overflow is not dependent on patient transfers; instead, it helps to contextualize other metrics. A non-zero system-wide overflow indicates that the system as a whole is overburdened and therefore desirable solutions are harder to find, and it represents a lower bound on the total overflow, although in general this bound is not achievable because we only consider transfers of newly admitted patients.

Table 1: General information on the models

Patient Modeling Assumptions	<ul style="list-style-type: none"> <li>• We consider only new patients as potential transfers between nodes.</li> <li>• We assume full knowledge about the number of new patients visiting a hospital each day.</li> <li>• We assume that the only resources limiting patient care are hospital beds, which are fixed in number at each node.</li> <li>• We assume that a fixed proportion of hospital beds are available to COVID-19 patients.</li> <li>• All patients have a length of stay governed by a distribution <math>\mathcal{L}_g</math>.</li> </ul>
Nurse Modeling Assumptions	<ul style="list-style-type: none"> <li>• Nurses may be moved between hospitals.</li> <li>• We have full knowledge of the initial number of nurses in each region, which is constant except for our reallocation.</li> <li>• In addition to bed availability, nurse availability limits patient care.</li> <li>• Some fixed proportion of nurses are available to treat COVID-19 patients.</li> </ul>
Sets	<ul style="list-style-type: none"> <li>• <math>\mathcal{N}</math>: patient treatment nodes, indexed by <math>n \in \mathcal{N}</math></li> <li>• <math>\mathcal{T}</math>: modeling days, indexed by <math>t \in \mathcal{T} = \{1, 2, 3, \dots, T\}</math></li> <li>• <math>\mathcal{G}</math>: patients groups, indexed by <math>g \in \mathcal{G}</math></li> <li>• <math>\mathcal{B}</math>: bed types, indexed by <math>\beta \in \mathcal{B}</math></li> </ul>

## Data

- $p_{g,i,t}$ : number of patients admitted to node  $i \in \mathcal{N}$  at time  $t \in \mathcal{T}$ , specified in group  $g \in \mathcal{G}$ . (zero for  $g$ s where  $\text{indegree}(g) \neq 0$ )
  - $p_{g,i,0}$ : number of initial patients in group  $g \in \mathcal{G}$  at node  $i \in \mathcal{N}$  at time  $t = 0$
  - $d_{g,i,t}$ : number of initial patients who were discharged from group  $g \in \mathcal{G}$  at node  $i \in \mathcal{N}$  at time  $t \in \mathcal{T}$
  - $b_{\beta,i}$ : number of beds of type  $\beta \in \mathcal{B}$  available for COVID-19 patients at node  $i \in \mathcal{N}$
  - $n_i$ : initial number of nurses at node  $i \in \mathcal{N}$
  - $G$ : directed graph where  $V(G) = \mathcal{N}$  and  $(i, j) \in E(G)$  if and only if node  $i$  may transfer resources to node  $j$
  - $\ell_g$ : distribution over length of stay for patients in group  $g \in \mathcal{G}$
  - $\mathcal{L}_g$ : cumulative distribution over length of stay for patients in group  $g \in \mathcal{G}$
- 

## Parameters

- $R_{\text{thresh}}$ : load ratio at which balancing penalization begins
  - $S_{\text{min}}$ : minimum number of patients that can be included in a transfer
  - $T_{\text{switch}}$ : minimum number of days a node must wait between sending and receiving patients
  - $C_{\text{balance}}$ : cost coefficient for load-balancing penalty  $\phi_{i,t}$
  - $C_{\text{smooth}}$ : cost coefficient for the smoothing penalty  $\delta_{i,j,t}$
  - $C_{\text{sent}}$ : cost coefficient for patient transfers  $s_{i,j,t}$
  - $C_{\text{setup}}$ : cost coefficient for the transfer indicator  $\rho_{i,j}$
  - $C_{\text{patient}}$ : cost coefficient for patient overflow
  - $C_{\text{nurse}}$ : cost coefficient for nurse overflow
  - $G_{\text{group}}$ : a graph where for all  $\forall i, j \in \mathcal{G}$ ,  $i \sim j$  if and only if patients from group  $i \in \mathcal{G}$  are transferred to group  $j \in \mathcal{G}$
  - $f$ : function mapping  $\mathcal{G}$  to  $\mathcal{B}$
  - $Q_{\beta}$ : ratio of nurse-days to patient-days for bed type  $\beta \in \mathcal{B}$
-

## Variables

- $s_{g,i,j,t}$ : number of patients of patient group  $g \in \mathcal{G}$  sent from node  $i \in \mathcal{N}$  to node  $j \in \mathcal{N}$  at time  $t \in \mathcal{T}$
  - $o_{\beta,i,t}$ : dummy variable for patient overflow in bed type  $\beta \in \mathcal{B}$  at node  $i \in \mathcal{N}$  at time  $t \in \mathcal{T}$
  - $\delta_{i,j,t}$ : dummy variable for the absolute difference in the number of patients sent from node  $i \in \mathcal{N}$  to node  $j \in \mathcal{N}$  between days  $t - 1$  and  $t \in \mathcal{T} \setminus \{1\}$ .
  - $\phi_{\beta,i,t}$ : dummy variable for the amount by which patient load ratio exceeds  $R_{\text{load}}$  at node  $i \in \mathcal{N}$  at time  $t \in \mathcal{T}$
  - $\rho_{i,j}$ : binary dummy variable which is equal to 1 if and only if there is a patient transfer between node  $i \in \mathcal{N}$  and node  $j \in \mathcal{N}$
  - $\nu_{1,i,t}, \nu_{2,i,t}$ : binary dummy variables used to enforce the minimum number of days a node must wait between sending and receiving patients
  - $\sigma_{i,j,t}$ : variable for nurses sent from region  $i \in \mathcal{N}$  to region  $j \in \mathcal{N}$  at time  $t \in \mathcal{T}$
  - $\theta_{i,t}$ : dummy variable for nurse overflow in region  $i \in \mathcal{N}$  at time  $t \in \mathcal{T}$
- 

## Expressions

- $\alpha_{g,i,t}$ : expression for the total number of active patients in group  $g \in \mathcal{G}$  at node  $i \in \mathcal{N}$  at time  $t \in \mathcal{T}$
  - $\chi_{g,i,t}$ : expression for the total number of patients entering group  $g \in \mathcal{G}$  at node  $i \in \mathcal{N}$  at time  $t \in \mathcal{T}$
  - $\gamma_{g,i,t}$ : expression for the total number of patients leaving group  $g \in \mathcal{G}$  at node  $i \in \mathcal{N}$  at time  $t \in \mathcal{T}$
  - $\eta_{i,t}$ : expression for the total number of nurses at node  $i \in \mathcal{N}$  at time  $t \in \mathcal{T}$
  - $q_{i,t}$ : expression for the total nurse demand at node  $i \in \mathcal{N}$  at time  $t \in \mathcal{T}$
- 

Optimal objective function value is also evaluated, which takes into account the overflow as well as the values of the logistical penalties associated with the solution. For each node-day we also consider the patient load, defined as the number of active patients divided by the number of beds. This load is therefore a measure of the stress on a node that is normalized to its capacity, and is therefore directly comparable between node-day pairs. The maximum load over time is an important metric as it measures how bad the stress on a node will get at the peak, which is what that node will have to prepare for. Finally, we evaluate solutions on the size of the patient transfers they make as well. Solutions that transfer fewer patients are more efficient and easier to implement in practice. As with overflow, the size of individual transfers matter, so we consider mean, median, and maximum non-zero transfer size. While these metrics capture many important aspects of the performance of a solution, there are other operational considerations, such as robustness, that are not fully captured. These secondary solution characteristics can be measured using additional metrics, and be qualitatively evaluated using the plots included in Section 5.



## 4 Data

The models developed in Section 3 require the following inputs in order to solve the problem of optimal COVID-19 patient redistribution:

1. The number of active COVID patients at each location in the study on the first day of the study period.
2. The proportion of these initial patients that were discharged at each location on each day.
3. The number of COVID patients admitted to each location on each day, which can optionally be uncertain within some certain interval.
4. The non-surge capacity of each location available for COVID patients.
5. The distribution over the length of stay for COVID patients.
6. A graph representing which pairs of locations are permitted to transfer patients with each other, which can optionally be directed.

In addition to such data, the group models require the capacity per bed type (ICU or ward) and the COVID patient inputs per patient group. Models involving nurse allocation also require nurse supply and demand, which is computed from the number of nurses at each location, the number of hours a nurse works per week, and the number of nurse-hours that should be devoted to each COVID patient per day.

While hospitals in the United States are mandated to report detailed metrics about their COVID patient load and response to the federal government, this data is aggregated at the state level before it is shared publicly, and even then only a small subset of the collected metrics are reported. Beyond this, there is no standardized publicly-available reporting of hospitalization data associated with COVID-19. However, some state and local governments have elected to release some of this data, enabling us to consider these regions as more adequate case studies for the methodology proposed in this work. Due to the lack of standardization in data reporting and the stringent data requirements of our models, data collection, processing, and cleaning proved to be a critical and substantial task. To aid in future studies regarding hospitalizations associated with COVID-19, we have open-sourced the data we collected and the code to process it<sup>2</sup>, as well as published a website that compiles relevant data sources<sup>3</sup>.

In this work, we present case studies based on historical data. In practice, however, this method would need to be implemented prospectively, using projections of the number of COVID-related hospitalizations. Forecasting COVID-related hospitalizations is beyond the scope of this work, but there exist multiple methods for such forecasting at both the hospital level and the state level (CDC 2020d, Predictive-Healthcare 2020). The reader is referred to such models for more information.

Table 2: Data sources used throughout this work.

Description	Source	Source URL
Florida Hospitalizations	Florida AHCA	<a href="https://ahca.myflorida.com/covid-19_alerts.shtml">https://ahca.myflorida.com/covid-19_alerts.shtml</a>
New Jersey Hospitalizations	New Jersey DOH	<a href="https://services7t.arcgis.com/Z0rixL1ManVefxqY/arcgis/rest/services/PPE_Capacity/FeatureServer/0">https://services7t.arcgis.com/Z0rixL1ManVefxqY/arcgis/rest/services/PPE_Capacity/FeatureServer/0</a>
Texas Hospitalizations	Texas DOH	<a href="https://dshs.texas.gov/coronavirus/additionaldata/">https://dshs.texas.gov/coronavirus/additionaldata/</a>
Hospital Beds	Definitive Healthcare	<a href="https://coronavirus-resources.esri.com/datasets/1044bb19da8d4dbfb6a96eb1b4ebf629_0">https://coronavirus-resources.esri.com/datasets/1044bb19da8d4dbfb6a96eb1b4ebf629_0</a>
COVID LOS Distribution	Lewnard et al. (2020)	<a href="https://www.bmj.com/content/369/bmj.m1923">https://www.bmj.com/content/369/bmj.m1923</a>

<sup>2</sup><https://github.com/flixpar/covid-resource-allocation>

<sup>3</sup><https://jhu-covid-optimization.github.io/covid-data/>

#### 4.1 COVID-Related Hospitalizations

The primary input data for each of the models developed in this work includes the number of initial active COVID patients, the proportion of initial patients discharged each day, and the number of newly admitted COVID patients each day, for each location. While such data may be known or projected for a specific case study, usually, the number of active COVID patients is more readily available. Therefore, we estimate the number of admitted patients and the proportion of initial patients discharged on each day from the reported number of active COVID patients. Specifically, for each location, given a candidate solution for the number of admitted patients and the proportion of initial patients discharged on each day, the number of active patients on each day is computed using constraint (1d) assuming zero patients are transferred. The objective function is then the  $\ell_2$ -distance between the computed number of active patients on each day and the reported number of active patients, which is approximately minimized over the candidate admitted and discharged solutions using random search optimization.

For the specific case studies considered in this work, different sources of data are considered to obtain COVID-related hospitalization data (see Figure 2). The first case study we consider is New Jersey. We obtained the number of active COVID patients from the New Jersey Department of Health. This dataset also reports the number of COVID patient admissions, discharges, and bed availability by hospital. However, the reported number of admissions and discharges are not entirely consistent with the reported number of active COVID patients, and therefore we decided to estimate these quantities from the reported active COVID patients using the above method. The second case study we consider is Florida. The Florida Agency for Healthcare Administration reports the number of active patients by hospital, but only reports the number of active COVID patients by county. We therefore estimate the number of active COVID patients at each hospital by multiplying the total number of patients at that hospital by the proportion of all patients that have COVID in the county where the hospital is located. We have also found that the number of active patients in Florida has sporadic single-day outliers which appear to be reporting errors. To resolve this problem, we run a simple outlier detection and correction algorithm. For each day and location we compute the median and median absolute deviation (MAD) of active patients over the surrounding five-day period. If the number of active patients deviates from the median by more than ten times the MAD then we consider it to be an outlier, and replace it with the median. The final case study we consider is Texas. The Texas Department of Health reports the number of active COVID patients by Trauma Service Area (TSA). This data does not appear to require correction, but it does not include patient admissions, so we estimate the COVID patient admissions and discharges using the method above.

#### 4.2 Hospital Capacity

To determine the capacity of each location in our case studies to care for COVID patients we used the Definitive Healthcare USA Hospital Beds dataset (see Table 2) to get the number of ICU and ward beds at each hospital in the United States. This data was collected before the onset of COVID in the United States, so it represents the ordinary hospital capacity rather than the surge capacity. We assumed that 35% of ward beds and 50% of ICU beds could be made available for COVID patients. These parameters are estimated from the state-level COVID hospitalization data reported by the Center for Disease Control (CDC). In the patient allocation models we assume that the number of beds that can be made available for COVID patients determines the capacity. In the nurse allocation model we instead assume that the number of nurses is the limiting factor for capacity instead.

#### 4.3 Patient Length of Stay

The distribution over the length of stay (LOS) in ward and ICU beds for COVID patients is estimated by Lewnard et al. (2020). We use a Weibull( $\lambda = 12.88, k = 1.38$ ) distribution for ward patients and a Weibull( $\lambda = 13.32, k = 1.58$ ) for ICU patients. See Figure 1. We discretize these distributions for use in the model. In the group model, it is required that patients stay for two days in a non-ICU bed before ICU admittance (Wunsch et al. 2011) and five days after ICU discharge. (Tiruvoipati et al. 2017)

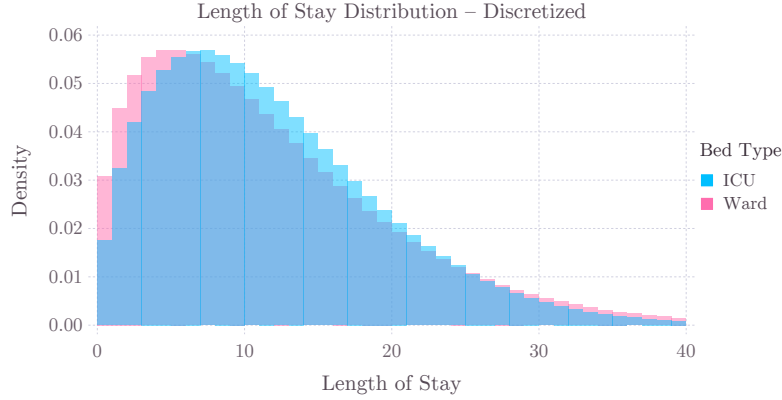


Figure 1: Distribution over the length of stay (LOS) for COVID-19 patients by bed type, discretized by day.

#### 4.4 Adjacency Graph

Our models constrain the solution so that patients are transferred only along the edges of a graph, which is an input to the model. In practice, we use this constraint to set an upper limit on the distance that a patient can be transferred by adding an edge between two locations in the graph if they are within a certain distance of each other. Given the latitude and longitude of the centroid of each location, we compute the distance matrix using the Haversine formula. We then threshold this distance matrix to determine the binary adjacency matrix for the patient transfer graph.

#### 4.5 Nurse Availability

In our New Jersey case study we validate our nurse and combined nurse plus patient models, which require data on the number of nurses that work at each hospital as well as the nurse demand. We estimate the number of nurses at each hospital from the New Jersey Department of Health’s Hospital Patient Care Staffing Report, which reports the mean number of patients per nurse. Multiplying this quantity by the mean number of patients-days per week we get the mean number of nurses-days per week. Assuming nurses work 36 hours per week on average, we multiply by 24/36 to get the number of total nurses. We estimate nurse demand from the number of active patients. The number of active patients times the number of nurses per patient yields the number of nurse-days required. We then compute the number of nurses required to supply that many nurse days per actual day, assuming that nurses work 36 hours per week on average.

In general, the number of nurses at each hospital would be known by the decision maker using this model, and therefore it would not have to be estimated this way.

#### 4.6 Data Collection Methods and Resources

In the process of collecting data for this work we found a great deal of data pertaining to COVID-19 or likely relevant to studies on COVID-19, from COVID-19 forecasts to nursing workforce surveys to health and hospital statistics. While there is a lot of available data for studies on COVID-19, there are also many gaps in the available data, which makes it difficult to know what data can be found and used without an extensive search. Much of the data is also not easily locatable as many sources show only a dashboard summarizing the available data, but the complete data can be accessed as well. We therefore have collected a list of over fifty data sources related to COVID-19 that we have used in this work or think are likely to be useful to other researchers. We also collect basic metadata about each data source, including a category, a description, an indication of whether the data is geographical, and the unit resolution. We are making this

list public on our website<sup>4</sup> so that other researchers may use it. We are also releasing the code used for this study, which contains scripts used to download, process, clean, and load the data we worked with<sup>5</sup>.

## 5 Results

In this section, we apply the models developed in Section 3 to a number of case studies in order to showcase the models' applicability to real examples, validate them, and demonstrate their effectiveness. In particular, we study the peak of the first wave of the COVID-19 pandemic in New Jersey and Texas. Miami-Dade county is also studied in the online supplement to this paper. We first evaluate each of the models we have developed to justify the inclusion of the additions we made to the base model. We then apply our final model to each of our case studies to demonstrate the potential effectiveness and impact of our approach.

### 5.1 Model Validation

To validate the models developed in this work, we first apply them to our New Jersey case study. New Jersey was selected for this task because the New Jersey Department of Health has published detailed data on COVID patients and bed availability at the hospital level. Specifically, we consider each of the 75 hospitals in the state of New Jersey from April 5th to June 15th, 2020. In Table 3 we evaluate each of the non-group patient allocation models according to the metrics defined in Section 3.4. We also include a "no transfer" model that does not transfer patients as a baseline against which to compare the effectiveness of our models. We note that without transfers, there is a significant overflow of COVID patients in patient-days, yet, as can be seen in Figure 6, the number of total active COVID patients is always less than the capacity of the whole system. Because of this we should expect to see a large reduction in the total overflow due to our models. Compared to the baseline, all models result in an overflow reduction of at least 28,862 patient-days, or 86.4%, which is a significant improvement for the hospital system.

The base model has the least restrictive constraints and penalties, and therefore is able to achieve the largest decrease in overflow, which is the primary goal. It also performs well according to other metrics including mean and median non-zero overflow. However, it does have some undesirable solution characteristics as compared to the other models, such as a higher percentage of hospital-days with an overflow, more total patients transferred, and a larger maximum-sized transfer. Figure 2 demonstrates that the base model often brings hospitals that were well under capacity all the way up to capacity as well, which in practice they will likely be unwilling to accommodate. The operational model addresses many of these issues, at the cost of greater total overflow. In particular, it transfers fewer patients overall and the mean, median, and maximum transfer size are all significantly smaller than those of the base model, which would make the transfers more operationally feasible.

Parameter and constraint selection is important for the operational model and their optimal choices are highly dependent on input data and model use case. In general, we believe that solution characteristics such as fewer transferred patients, no large spikes in the number of transferred patients, and maintaining a cushion between the number of active patients and the capacity are important because they make the solution more feasible in practice. Adding small penalties on the total amount of patients transferred and the smoothness in the number of patients transferred between days also seems to improve the solution quality while not significantly increasing the running time of the solver. Therefore,  $C_{\text{sent}} = 0.01$ ,  $C_{\text{smooth}} = 0.01$  are used in the operational model. We also include (3i) and (3j), which ensure that hospitals are not sent over or further over capacity by our model solution, as we believe this solution characteristic will be important to decision makers. On the other hand, setting a lower bound on the size of a patient transfer or adding a setup cost between hospitals may only make sense in some situations, and requires adding integer variables to the model which drastically increases the time it takes to solve. In our operational model we employ only the optional constraints and

<sup>4</sup><https://jhu-covid-optimization.github.io/covid-data/>

<sup>5</sup><https://github.com/flixpar/covid-resource-allocation>

penalties that have broad applicability, however, the other optional constraints and penalties remain in our model to provide flexibility to practitioners.

In Figure 2 we compare the number of active patients at a subset of the hospitals under the base model, the operational model, and without transfers. It can be seen in this figure how the models perform in utilizing excess capacity in the system to accommodate the surge in demand on regional nodes. In this example, both the models achieve their goal of eliminating the overflow in all the hospitals after a couple of days of redistributing the patients. The base model on occasion sends hospitals that were not experiencing an overflow over capacity, which in practice hospitals may not be willing to accommodate. The operational model on the other hand maintains a 5% buffer between the number of active patients and the capacity for each hospital-day.

The robust model also addresses some of the shortcomings of the base model, however, its primary purpose is to protect the feasibility of the solution against different prospective. In Figure 5 we see the results of the base, robust, and no-transfer models number with active patient counts sampled uniformly from the uncertainty set. In this figure we can see scenarios where the base model sends a hospital over capacity while the robust model does not because the base model only considers the nominal number of admitted patients. The safety associated with the robust model will be its main value to decision makers, who must account for uncertainty and minimize the negative impact of their actions and policies. It is also important to note that all models perform at least as well as the no-transfer model in all possible scenarios within the uncertainty set. Figure 4 plots the outcome in terms of total overflow from 400 scenarios sampled uniformly from the uncertainty set, and demonstrates that for all of these cases both the base model and the robust model perform better than the baseline. In this figure we also notice that the robust model consistently performs worse than the base model, which is an unavoidable result of robustness shrinking the the LP's feasible region. Fortunately, we see that this gap is relatively small compared with the gap between each model and the no-transfer model.

The group patient allocation models add the capability to differentiate between patients in different care-paths and to have multiple bed types. We ran the group patient allocation models on the New Jersey hospital-level data at a higher resolution, considering both ward and ICU patients, and in Figure 4 we compare the results for each group. These figures show that the base group model is able to achieve a smaller overflow yet has to transfer more patients in both the ICU and ward blocks than either the no-transfer model or the group operational model. Interestingly, we see that it is not possible to reduce overflow as much in the ward as in the ICU. This is because all patients that visit the ICU also stay in the ward for seven days (two before ICU, five after ICU) but can only be transferred at the beginning of their stay, which means that the number of patient-days in the ward is much higher and there is less flexibility to transfer those extra patient-days.

The last set of models we analyze here are the nurse allocation models. There is a large shortage of nurses-days in the case study we consider because of the increased patient load. Since nurses are critical to properly care for all patients, represents a huge issue. A nurse shortage can be alleviated by increasing the number of hours that each nurse works, utilizing part-time nurses more, or decreasing the care given to patients, however, these solutions are clearly undesirable because they place a larger burden on nurses and have the potential to hurt patient outcomes. Transferring nurses is an appealing alternative to these choices. We see in Table 5 that our nurse allocation model is able to alleviate much of the nurse shortage.

## 5.2 Patient Redistribution Case Studies

In this section we apply the operational model to our two primary case studies, New Jersey and Texas. Compared to the previous section, in which our focus was on validation of our approach and models, here we focus on the solution and operational insights.

Table 3: Evaluation of the performance of each of the non-group patient allocation models in our New Jersey case study.

	No Transfers	Base Model	Operational Model	Base Robust Model	Operational Robust Model
Overflow	33406	3800	3930	4173	4544
Overflow Reduction	0.0%	88.62%	88.24%	87.51%	86.4%
Median Non-Zero Overflow	27.5	8	11	21	9
Mean Non-Zero Overflow	44.7	27.7	30.2	35.4	26.6
Max Non-Zero Overflow	244	220	220	220	220
Median Load	40.98%	51.8%	52.0%	53.0%	53.52%
Mean Load	62.15%	55.2%	54.91%	55.18%	56.54%
Max Load	2000.0%	400.0%	400.0%	400.0%	900.0%
Percent Of Hospital-Days With An Overflow	13.85%	2.54%	2.41%	2.19%	3.17%
Total Patients Transferred	0	4298	3932	4894	4255
Percent Of Patients Transferred	0.0%	14.25%	13.04%	16.23%	14.11%
Median Non-Zero Transfer	0	7	1	5	2
Mean Non-Zero Transfer	0	9.7	2.4	7.8	2.6
Max Non-Zero Transfer	0	55	19	52	25
Percent Of Hospital-Days With A Transfer	0.0%	12.81%	24.31%	17.35%	24.65%

Table 4: Evaluation of the performance of each of the group patient allocation models in our New Jersey case study.

	ICU			Ward		
	No Transfers	Base Group Model	Operational Group Model	No Transfers	Base Group Model	Operational Group Model
Overflow	13452	1129	3189	37068	18270	28912
Overflow Reduction	0.0%	91.6%	76.3%	0.0%	50.7%	22.0%
Median Non-Zero Overflow	8	1	2	25	16	26
Mean Non-Zero Overflow	12	3.6	4.1	37.3	24	37.3
Max Non-Zero Overflow	89	37	37	231	148	226
Median Load	50.0%	78.6%	75.0%	62.7%	73.3%	68.7%
Mean Load	84.4%	71.9%	73.2%	82.0%	76.6%	79.5%
Max Load	735.7%	433.3%	433.3%	591.5%	414.9%	580.9%
Percent Of Hospital-Days With An Overflow	32.5%	9.1%	22.6%	28.8%	22.0%	22.5%
Total Patients Transferred	0	1432	1016	0	3060	1395
Percent Of Patients Transferred	0.0%	39.4%	27.9%	0.0%	17.6%	8.0%
Median Non-Zero Transfer	0	1	1	0	3	1
Mean Non-Zero Transfer	0	2	1.4	0	4.8	2.2
Max Non-Zero Transfer	0	13	13	0	52	36
Percent Of Hospital-Days With A Transfer	0.0%	29.9%	25.2%	0.0%	24.3%	18.9%

### 5.2.1 Case Study 1: New Jersey.

In the previous section we used the New Jersey case study to compare models, here we investigate the solution found by the operational model and its significance for New Jersey. Once again we consider each of the 75 hospitals in the state of New Jersey from April 5th to June 15th, 2020.

As noted before, a number of the New Jersey hospitals experienced a COVID patient overflow and as a result, the hospitals took action to increase their original patient capacity to accommodate the additional patients. We estimate that they had to add at least 33,406 bed-days of capacity, or at least a total of 2093 beds. According to a report from the New Jersey Hospital Association (New Jersey Hospital Association 2020), hospitals actually added at least 2800 ICU beds. Even with the extra capacity, the report also states that the overall load on the system increased from its typical value of



Figure 2: Active patients by model for a representative sample of nine hospitals in New Jersey.

Table 5: Evaluation of the performance of each of the nurse allocation model in our New Jersey case study.

	No Transfers	Base Nurses Model	Operational Nurse Model
Shortage	55119	34354	34629
Shortage Reduction	0.0%	37.7%	37.2%
Median Non-Zero Shortage	38	22	22
Mean Non-Zero Shortage	47	32	31.9
Max Non-Zero Shortage	145	202	203
Percent Of Hospital-Days With A Shortage	67.9%	62.0%	62.8%
Total Nurse Transfers	0	8851	1025
Median Non-Zero Transfer	0	8	1
Mean Non-Zero Transfer	0	18.2	3.1
Max Non-Zero Transfer	0	195	74
Percent Of Hospital-Days With A Transfer	0.0%	32.6%	22.4%

62% to 82%, meaning that there was also an increase in demand for nurses. Evidently, the New Jersey hospital system faced an extreme pressure during the first wave of the pandemic.

Despite this extreme pressure, our model found a solution that would have eliminated more than 88% of the total overflow, as can be seen in Table 7. This translates to 3930 surge bed-days and just 1141 additional beds, which is a huge improvement over not transferring patients. It is clear from these metrics that New Jersey hospitals would not have been forced to create nearly as much surge capacity, putting a far smaller strain on the system, had they

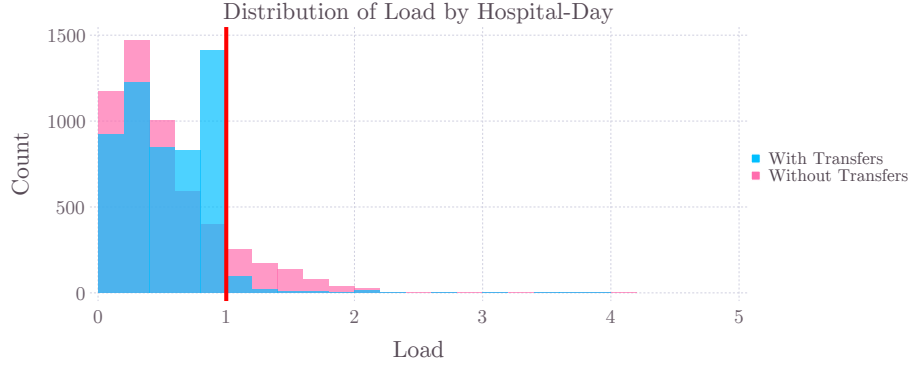


Figure 3: Distribution over the COVID patient load by hospital-day in New Jersey with the operational model.

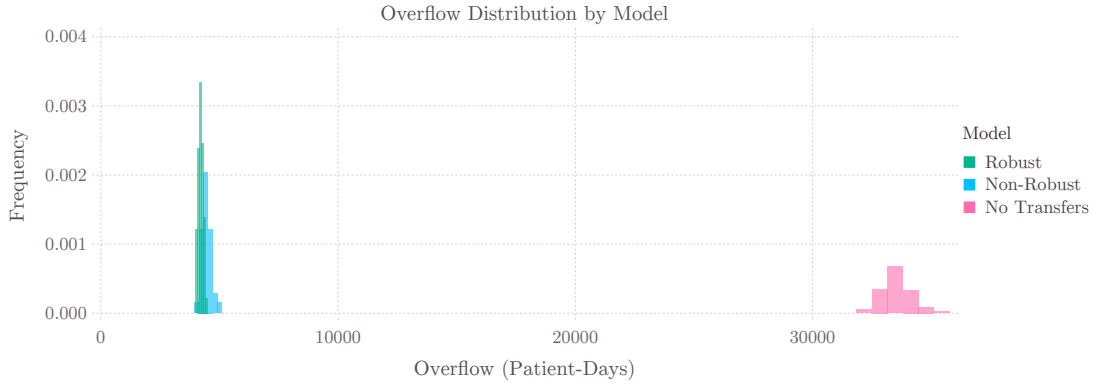


Figure 4: Distribution over the total overflow for all hospitals in New Jersey for three models under scenarios sampled uniformly from the uncertainty set.

made optimal patient transfers. This solution would have involved transferring 13% of all COVID patients in New Jersey. While this represents a very large number of people, we believe it is an acceptable trade-off for the potential improvements in patient care, system efficiency, and lower costs to create surge capacity. Clearly there is a trade-off between the number of transfers and the amount of overflow reduction, which can be made by increasing the penalty on the number of transferred patients in the model. We selected to use a small penalty just to encourage the model not to make unnecessary transfers, but this is potentially an important tool for decision makers to find operationally feasible solutions.

Figure 7 plots the COVID patient load over time for a representative sample of hospitals, demonstrating that the model is able to get the additional demand under control almost immediately, and is able to maintain a 5% buffer between demand and capacity for most hospital-days. In Figure 8 we see that the overflow is primarily clustered in the northeast, close to New York City which was the epicenter of the pandemic during this time period, and that the hospitals in the rest of the state were able to stay under capacity. It is this imbalance, and the resulting transfers out of the northeast region, that made the model so successful.

### 5.2.2 Case Study 2: Texas.

The second case study considered in this section is the state of Texas at the Trauma Service Area (TSA) level from June to August 2020 when regions of Texas experienced a severe wave of increased COVID-19 cases and some hospitals





Figure 5: Active COVID patients over time under five scenarios sampled uniformly from the robust uncertainty set at a representative sample of nine hospitals in New Jersey.

Table 6: Parameters for each case study.

Region	New Jersey	Florida	Texas
Allocation level	Hospital	Hospital	TSA
Bed type	All	All	All
Start date	2020-04-05	2020-07-04	2020-06-15
End date	2020-06-15	2020-08-10	2020-08-15
Number of days	72	38	62
Number of hospitals	75	24	22
Number of hospital-days	5400	912	1364

even practiced transferring patients to balance their load. TSAs are the smallest level at which the state of Texas reports hospitalizations. Analyzing Texas at the TSA level demonstrates that our methods can be applied to healthcare systems at a coarser level while remaining effective and valuable. This approach may be valuable in cases where data is unavailable at more local levels or when entire regions of hospitals are all out or nearly out of capacity. However, such an approach implicitly assumes that there is in fact an optimal distribution of patients among hospitals in each TSA. While this assumption does not hold in general, decision makers can use the model for each region with particularly unbalanced COVID patient load given the results of the system-wide model.

Figure 11 shows that a number of the TSAs go well over a patient load ratio of 1.0, meaning that they had to create a significant amount of surge capacity to care for the COVID patients. However, it should be noted that, similar to the case of the state of new jersey, the surges in the demand were regional and the state did not surpass its capacity as a whole. This, in turn, enabled the models to reduce the total patient overflow to zero by transferring patients among nodes. The results of the transfers can be seen in Table 7 and Figure 11. The operational model was also able to maintain a gap of

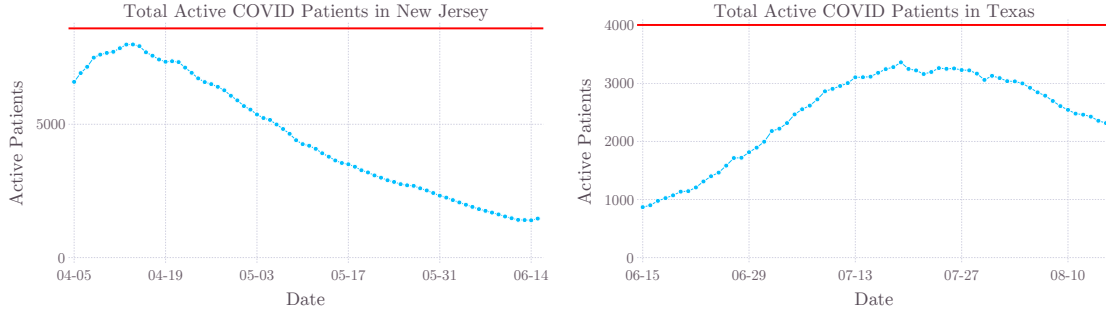


Figure 6: Total active COVID patients (blue) versus the bed capacity for COVID patients (red) for New Jersey and Texas.

Table 7: Evaluation of the performance of the operational model in each case study.

Case Study	New Jersey	Texas
Overflow	3930	21
Overflow Reduction	88.24%	99.84%
Number Of Hospital-Days With An Overflow	130	10
Percent Of Hospital-Days With An Overflow	2.41%	0.73%
Total Patients Transferred	3932	1705
Percent Of Patients Transferred	13.04%	12.1%
Median Non-Zero Transfer	1	2
Mean Non-Zero Transfer	2.4	3.5
Max Non-Zero Transfer	19	23
Percent Of Hospital-Days With A Transfer	24.31%	39.81%

5% of capacity between the number of active patients and the capacity. The solution of this particular model involved transferring 1765 patients, or 12.1% of the total COVID patient population for this time period.

### 5.3 Case Study 3: Miami

Our final case study considers hospitals in Florida in July and early August 2020, encompassing the peak of the first wave in the state. We specifically target the 24 Class I hospitals in Miami-Dade County as it was among the most severely impacted counties in Florida at the time. We therefore are investigating the potential value of performing patient re-distribution in a local hospital system rather than state-wide. We ran the operation model on this system with the same parameters as in the other case studies. The results of this model can be seen in Table 8.

According to the evaluation metrics reported in Table 8, the operational model could have reduced COVID patient overflow in Miami-Dade by nearly 90% while requiring fewer than 4% of the total COVID patients to be transferred. Such a large decrease in overflow and small number of required transfers represents an easy and effective way to accommodate the surge of COVID patients that the system had to face. The data shows that these hospitals had to increase their capacity by 1251 patients at the peak of the pandemic to accommodate COVID patients. Being able to reduce this by as much as 90% would have significantly reduced the burden on many hospitals and made ensuring that all patients were properly cared for much easier, cheaper, and more efficient.

### 5.4 Implementation

All models described in Section 3 were implemented in Julia 1.5 using the JuMP library for modeling (Dunning et al. 2017). Models were solved using Gurobi 9.0.3 with default options. The code is publicly available<sup>6</sup>.

<sup>6</sup><https://github.com/flixpar/covid-resource-allocation>

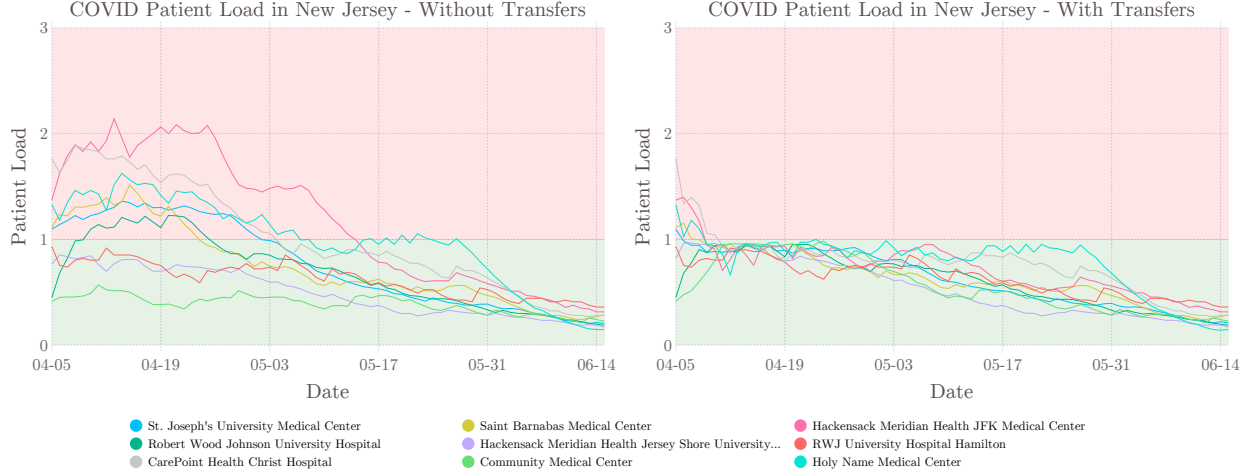


Figure 7: COVID patient load at a representative sample of 9 hospitals in New Jersey, with and without transfers.

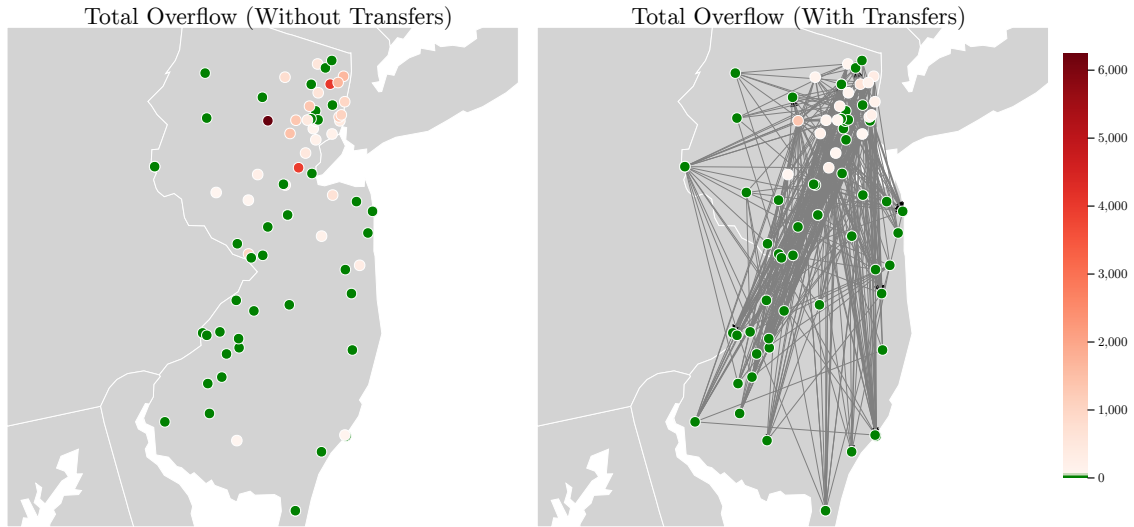


Figure 8: COVID patient overflow by hospital in New Jersey, with and without transfers. Edges indicate patient transfers.

## 5.5 Website

We have demonstrated the potential of our methods to reduce the burden on hospitals facing an extreme surge in demand. In order to make this a practical tool that hospital systems or governments can use we have publicly released the code and data that we have used. In addition to this, we have deployed an interactive website that people can use to explore the potential impact that optimal patient transfers could have. It allows users to specify the region, time period, and model parameters, and run our model in real time. It then displays metrics and animated figures which enables users to effectively understand and evaluate the solution without detailed knowledge of the workings of the model. The website may be found at: <https://covid-hospital-operations.com/>.



Figure 9: Nurse supply and demand for COVID-19 patients at a representative subset of nine hospitals in New Jersey.

Table 8: Evaluation of the performance of the operational model in the Florida case study.

Overflow	51
Overflow Reduction	89.7%
Median Non-Zero Overflow	1
Mean Non-Zero Overflow	1.7
Max Non-Zero Overflow	21
Median Load	88.0%
Mean Load	85.8%
Max Load	121.0%
Percent Of Hospital-Days With An Overflow	3.3%
Total Patients Transferred	251
Percent Of Patients Transferred	3.9%
Median Non-Zero Transfer	1
Mean Non-Zero Transfer	0.6
Max Non-Zero Transfer	6
Percent Of Hospital-Days With A Transfer	42.1%

## 6 Conclusion

In this work, we introduced a methodology that can be used to redistribute demand and resources among hospitals during high demand periods such as the COVID-19 pandemic. The results of applying these methods to real examples from the first wave of the pandemic showed that redistribution effectively helps hospital systems balance patient loads and potentially provide better care, while reducing the amount of surge capacity required. Additionally, such a framework promotes systematic collaboration across healthcare entities to proactively respond to predicted high demand events.

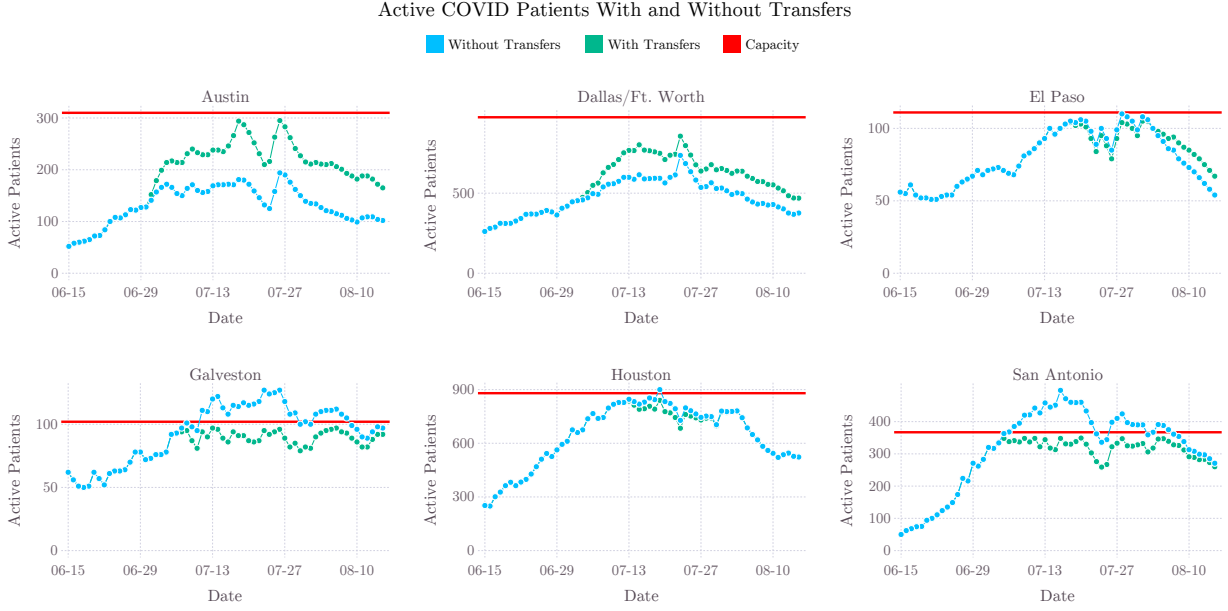


Figure 10: Active COVID patients for a subset of TSAs in Texas, with and without patient transfers.

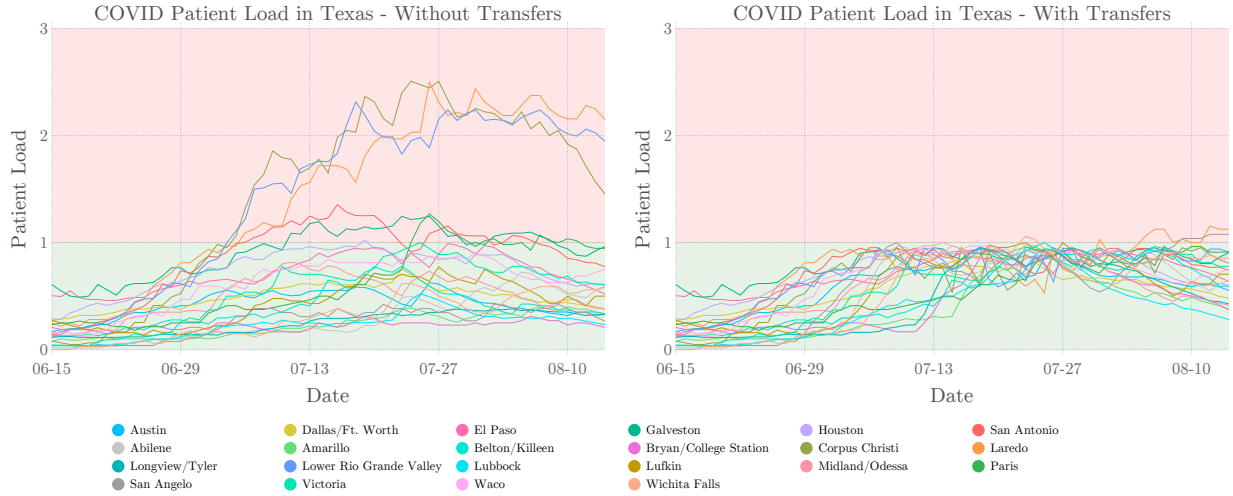


Figure 11: COVID patient load by TSA, with and without transfers.

The problem of optimizing hospital system decisions during a stressful time like a pandemic is complex and needs proper attention to operational constraints specific to this particular problem. As such, in this work we introduced a flexible model that can accommodate the varying requirements of hospital systems and optimally allocate newly admitted patients among them to balance patient loads and consequently, promote patient care quality while reducing operating costs. The results of applying these models retrospectively to real world systems at different levels of analysis demonstrate successful load balancing, overflow reduction, and good operational characteristics. In order to improve the practical utility of this work, we also considered resource constraints (specifically the number of available nurses), included different care paths into the model, and made the model robust against some sources of uncertainty. Even incorporating all such additions and constraints, the models outperform no redistribution of patients by a large margin.

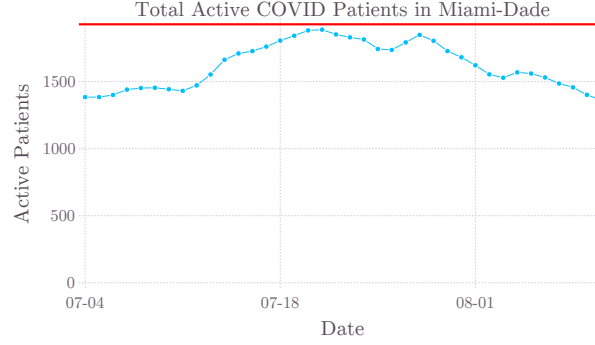


Figure 12: Total active COVID patients (blue) versus the bed capacity for COVID patients (red) in Miami-Dade County.

Although all the models perform exceptionally in reducing the burden on entities of a healthcare system, there are some limitations to the models that must be addressed. First, primarily using LP formulations (and some MILP) enables us to find solutions efficiently for large problems, but we are unable to accomplish quadratic load balancing or individual patient tracking using such schemes. Additionally, the scarcity and inherent uncertainty in COVID-related data also poses challenges to using such models. Specifically, due to the decentralized response of the United States to COVID-19, the quality of infection, hospitalization, personnel, and resource data varies widely among counties, states, and regions, but is often poor. Data quality is a serious issue because it is difficult to take decisive action without full information about the situation, although this issue can be alleviated somewhat using our robust model. It should also be noted that the operational constraints provided do not and can not capture the full range of operational considerations that decision makers will need to consider. However, we have formulated the models in such a way that extensions are simple.

The primary takeaway from this work should be that patient and resource re-distribution can be a very effective tool to reduce the burden on hospitals under stress from COVID-19, and that the models we present can solve this model optimally and in an operationally feasible manner.

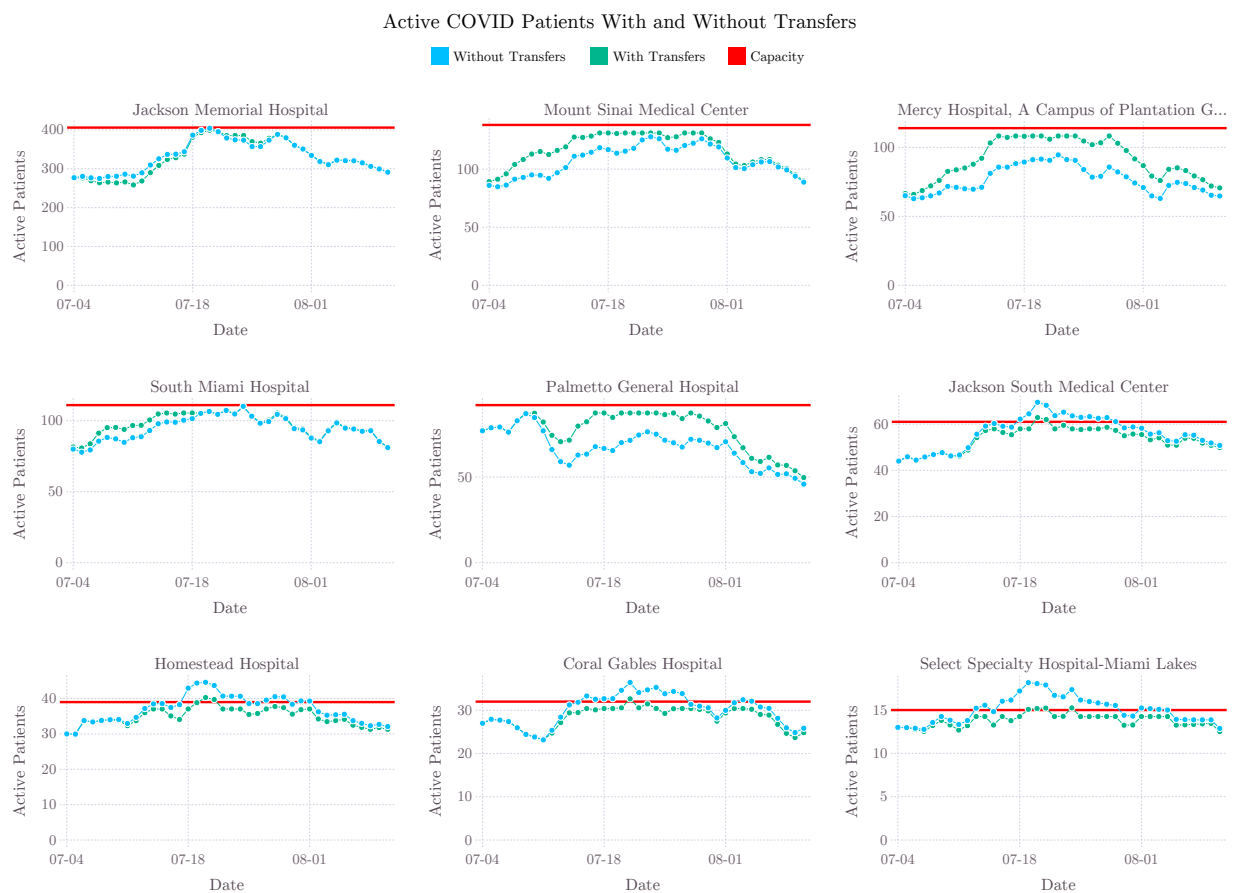


Figure 13: Active COVID patients for 9 representative hospitals in Miami-Dade County, with and without patient transfers.

## References

- David S Hui, Esam I Azhar, Tariq A Madani, Francine Ntoumi, Richard Kock, Osman Dar, Giuseppe Ippolito, Timothy D Mchugh, Ziad A Memish, Christian Drosten, et al. The continuing 2019-ncov epidemic threat of novel coronaviruses to global health—the latest 2019 novel coronavirus outbreak in wuhan, china. *International Journal of Infectious Diseases*, 91:264–266, 2020.
- Ensheng Dong, Hongru Du, and Lauren Gardner. An interactive web-based dashboard to track covid-19 in real time. *The Lancet infectious diseases*, 20(5):533–534, 2020.
- IHME, Christopher JL Murray, et al. Forecasting covid-19 impact on hospital bed-days, icu-days, ventilator-days and deaths by us state in the next 4 months. *MedRxiv*, 2020.
- Sarah Mervosh, Denise Lu, and Vanessa Swales. See which states and cities have told residents to stay at home. *New York Times*, 2020.
- CDC. Healthcare facilities: Managing operations during the covid-19 pandemic, Jun 2020a. URL <https://www.cdc.gov/coronavirus/2019-ncov/hcp/guidance-hcf.html>.
- CDC. Strategies for optimizing the supply of facemasks, 2020b.
- Marco Varkevisser, Stéphanie A van der Geest, and Frederik T Schut. Do patients choose hospitals with high quality ratings? empirical evidence from the market for angioplasty in the netherlands. *Journal of Health Economics*, 31(2):371–378, 2012.
- Florian Dreys. How patients choose hospitals: Using the stereotypic content model to model trustworthiness, warmth and competence. *Health services management research*, 26(2-3):95–101, 2013.
- Peter T Vanberkel, Richard J Boucherie, Erwin W Hans, Johann L Hurink, and Nelly Litvak. Efficiency evaluation for pooling resources in health care. *OR spectrum*, 34(2):371–390, 2012.
- Sarah Sims, Gillian Hewitt, and Ruth Harris. Evidence of collaboration, pooling of resources, learning and role blurring in interprofessional healthcare teams: a realist synthesis. *Journal of Interprofessional Care*, 29(1):20–25, 2015.
- Jiri Chod and Nils Rudi. Resource flexibility with responsive pricing. *Operations Research*, 53(3):532–548, 2005.
- Leon Boudourakis, David M. Silvestri, Shaw Natsui, R. James Salway, Mona Krouss, Amit Uppal, Alex Izaguirre, Mathew Siegler, Sonya Bernstein, Katelyn Prieskorn, Akshatta Dahake, and Eric K. Wei. Using interfacility transfers to 'level-load' demand from surging covid-19 patients: Lessons from nyc health hospitals, Jul 2020. URL <https://www.healthaffairs.org/doi/10.1377/hblog20200710.163676/full/>.
- CDC. Key considerations for transferring patients to relief healthcare facilities when responding to community transmission of covid-19 in the united states, 2020c. URL <https://www.cdc.gov/coronavirus/2019-ncov/hcp/relief-healthcare-facilities.html>.
- Tarek Hegazy. Optimization of resource allocation and leveling using genetic algorithms. *Journal of construction engineering and management*, 125(3):167–175, 1999.
- Jinke Ren, Guanding Yu, Yunlong Cai, and Yinghui He. Latency optimization for resource allocation in mobile-edge computation offloading. *IEEE Transactions on Wireless Communications*, 17(8):5506–5519, 2018.
- G Kuchuk, S Nechausov, and V Kharchenko. Two-stage optimization of resource allocation for hybrid cloud data store. In *2015 International Conference on Information and Digital Technologies*, pages 266–271. IEEE, 2015.
- Jose Cruz, Genshe Chen, Dongxu Li, and Xu Wang. Particle swarm optimization for resource allocation in uav cooperative control. In *AIAA Guid, Ivan ance, Navigation, and Control Conference and Exhibit*, page 5250, 2004.
- George Tychogiorgos and Kin K Leung. Optimization-based resource allocation in communication networks. *Computer Networks*, 66:32–45, 2014.
- Dimitris Bertsimas and Melvyn Sim. The price of robustness. *Operations research*, 52(1):35–53, 2004.
- Virginie Gabrel, Cécile Murat, and Aurélie Thiele. Recent advances in robust optimization: An overview. *European journal of operational research*, 235(3):471–483, 2014.
- Mehdi Najafi, Kourosh Eshghi, and Wout Dullaert. A multi-objective robust optimization model for logistics planning in the earthquake response phase. *Transportation Research Part E: Logistics and Transportation Review*, 49(1):217–249, 2013.
- Ruth Luscombe and Erhan Kozan. Dynamic resource allocation to improve emergency department efficiency in real time. *European Journal of Operational Research*, 255(2):593–603, 2016.
- Mojisola Otegbeye, Roslyn Scriber, Donna Ducoin, and Amy Glasofer. Designing a data-driven decision support tool for nurse scheduling in the emergency department: a case study of a southern new jersey emergency department. *Journal of emergency nursing*, 41(1):30–35, 2015.



- Nelly Litvak, Marleen Van Rijsbergen, Richard J Boucherie, and Mark van Houdenhoven. Managing the overflow of intensive care patients. *European journal of operational research*, 185(3):998–1010, 2008.
- Eric Toner and Richard Waldhorn. What hospitals should do to prepare for an influenza pandemic. *Biosecurity and Bioterrorism: Biodefense Strategy, Practice, and Science*, 4(4):397–402, 2006.
- John C Papageorgiou. Some operations research applications to problems of health care systems (a survey). *International journal of bio-medical computing*, 9(2):101–114, 1978.
- Abdur Rais and Ana Viana. Operations research in healthcare: a survey. *International transactions in operational research*, 18(1): 1–31, 2011.
- Margaret L Brandeau. Allocating resources to control infectious diseases. In *Operations Research and Health Care*, pages 443–464. Springer, 2005.
- Scott D Halpern and Franklin G Miller. The urge to build more intensive care unit beds and ventilators: Intuitive but errant, 2020.
- Ben S Cooper, Richard J Pitman, W John Edmunds, and Nigel J Gay. Delaying the international spread of pandemic influenza. *PLoS Med*, 3(6):e212, 2006.
- Ateev Mehrotra, Kristin Ray, Diane M Brockmeyer, Michael L Barnett, and Jessica Anne Bender. Rapidly converting to "virtual practices": outpatient care in the era of covid-19. *NEJM catalyst innovations in care delivery*, 1(2), 2020a.
- Ezekiel J Emanuel, Govind Persad, Ross Upshur, Beatriz Thome, Michael Parker, Aaron Glickman, Cathy Zhang, Connor Boyle, Maxwell Smith, and James P Phillips. Fair allocation of scarce medical resources in the time of covid-19, 2020.
- Douglas B White and Bernard Lo. A framework for rationing ventilators and critical care beds during the covid-19 pandemic. *Jama*, 323(18):1773–1774, 2020.
- Timothy J Judson, Anobel Y Odisho, Aaron B Neinstein, Jessica Chao, Aimee Williams, Christopher Miller, Tim Moriarty, Nathaniel Gleason, Gina Intinarelli, and Ralph Gonzales. Rapid design and implementation of an integrated patient self-triage and self-scheduling tool for covid-19. *Journal of the American Medical Informatics Association*, 27(6):860–866, 2020.
- Sen Pei, Sasikiran Kandula, Wan Yang, and Jeffrey Shaman. Forecasting the spatial transmission of influenza in the united states. *Proceedings of the National Academy of Sciences*, 115(11):2752–2757, 2018.
- Fotios Petropoulos and Spyros Makridakis. Forecasting the novel coronavirus covid-19. *PloS one*, 15(3):e0231236, 2020.
- Nicholas P Jewell, Joseph A Lewnard, and Britta L Jewell. Caution warranted: using the institute for health metrics and evaluation model for predicting the course of the covid-19 pandemic, 2020.
- Weston C Roda, Marie B Varughese, Donglin Han, and Michael Y Li. Why is it difficult to accurately predict the covid-19 epidemic? *Infectious Disease Modelling*, 2020.
- Matjaž Perc, Nina Gorišek Miksić, Mitja Slavinec, and Andraž Stožer. Forecasting covid-19. *Frontiers in Physics*, 8:127, 2020.
- Joseph A Lewnard, Vincent X Liu, Michael L Jackson, Mark A Schmidt, Britta L Jewell, Jean P Flores, Chris Jentz, Graham R Northrup, Ayesha Mahmud, Arthur L Reingold, et al. Incidence, clinical outcomes, and transmission dynamics of severe coronavirus disease 2019 in california and washington: prospective cohort study. *bmj*, 369, 2020.
- Eleanor M Rees, Emily S Nightingale, Yalda Jafari, Naomi R Waterlow, Samuel Clifford, Carl AB Pearson, Thibaut Jombart, Simon R Procter, Gwenan M Knight, CMMID Working Group, et al. Covid-19 length of hospital stay: a systematic review and data synthesis. 2020.
- Gary E Weissman, Andrew Crane-Droesch, Corey Chivers, ThaiBinh Luong, Asaf Hanish, Michael Z Levy, Jason Lubken, Michael Becker, Michael E Draugelis, George L Anesi, et al. Locally informed simulation to predict hospital capacity needs during the covid-19 pandemic. *Annals of internal medicine*, 2020.
- Hina Arora, TS Raghu, and Ajay Vinze. Resource allocation for demand surge mitigation during disaster response. *Decision Support Systems*, 50(1):304–315, 2010.
- Sanjay Mehrotra, Hamed Rahimian, Masoud Barah, Fengqiao Luo, and Karolina Schantz. A model of supply-chain decisions for resource sharing with an application to ventilator allocation to combat covid-19. *Naval Research Logistics*, 2020b.
- Ira M Longini Jr, Eugene Ackerman, and Lila R Elveback. An optimization model for influenza a epidemics. *Mathematical Biosciences*, 38(1-2):141–157, 1978.
- Lorenzo Lampariello and Simone Sagratella. Effectively managing diagnostic tests to monitor the covid-19 outbreak in italy. Technical report, Tech. rep. Optimization Online, 2020. url: <http://www.optimization-online...>, 2020.
- Mei Fong Liew, Wen Ting Siow, Ying Wei Yau, and Kay Choong See. Safe patient transport for covid-19. *Critical Care*, 24(1):1–3, 2020.
- Lihui Bai and Jiang Zhang. An incentive-based method for hospital capacity management in a pandemic: the assignment approach. *International Journal of Mathematics in Operational Research*, 6(4):452–473, 2014.

- Li Sun, Gail W DePuy, and Gerald W Evans. Multi-objective optimization models for patient allocation during a pandemic influenza outbreak. *Computers & Operations Research*, 51:350–359, 2014.
- Lucas Lacasa, Robert Challen, Ellen Brooks-Pollock, and Leon Danon. A flexible load sharing system optimising icu demand in the context of covid-19 pandemic. *medRxiv*, 2020.
- CDC. Covid-19 forecasts: Hospitalizations, September 2020d. URL <https://www.cdc.gov/coronavirus/2019-ncov/cases-updates/hospitalizations-forecasts.html>.
- Penn Medicine Predictive-Healthcare. Covid-19 hospital impact model for epidemics, March 2020. URL <https://code-for-philly.gitbook.io/chime/>.
- Hannah Wunsch, Derek C Angus, David A Harrison, Walter T Linde-Zwirble, and Kathryn M Rowan. Comparison of medical admissions to intensive care units in the united states and united kingdom. *American journal of respiratory and critical care medicine*, 183(12):1666–1673, 2011.
- Ravindranath Tiruvoipati, John Botha, Jason Fletcher, Himangsu Gangopadhyay, Mainak Majumdar, Sanjiv Vij, Eldho Paul, David Pilcher, Australia, and New Zealand Intensive Care Society (ANZICS) Clinical Trials Group. Intensive care discharge delay is associated with increased hospital length of stay: A multicentre prospective observational study. *PloS one*, 12(7):e0181827, 2017.
- New Jersey Hospital Association. The rise and fall of covid in new jersey: Hospitals respond as the curve flattens, May 2020. URL <http://www.njha.com/media/601210/Rise-and-Fall-of-COVID.pdf>.
- Iain Dunning, Joey Huchette, and Miles Lubin. Jump: A modeling language for mathematical optimization. *SIAM Review*, 59(2): 295–320, 2017.