

# On Fair Routing from Emergency Departments to Hospital Wards: QED Queues with Heterogeneous Servers

Avishai Mandelbaum

Faculty of Industrial Engineering and Management, Technion–Israel Institute of Technology, Haifa 32000, Israel,  
avim@tx.technion.ac.il

Petar Momčilović

Department of Industrial and Systems Engineering, University of Florida, Gainesville, Florida 32611,  
petar@ise.ufl.edu

Yulia Tseytlin

Faculty of Industrial Engineering and Management, Technion–Israel Institute of Technology, Haifa 32000, Israel; and  
IBM Research Lab, Haifa University, Haifa 31905, Israel, yuliatse@gmail.com

The interface between an emergency department and internal wards is often a hospital's bottleneck. Motivated by this interaction in an anonymous hospital, we analyze queueing systems with heterogeneous server pools, where the pools represent the wards, and the servers are beds. Our queueing system, with a single centralized queue and several server pools, forms an inverted- $V$  model. We introduce the randomized most-idle (RMI) routing policy and analyze it in the quality- and efficiency-driven regime, which is natural in our setting. The RMI policy results in the same server fairness (measured by idleness ratios) as the longest-idle-server-first (LISF) policy, which is commonly used in call centers and considered fair. However, the RMI policy utilizes only the information on the number of idle servers in different pools, whereas the LISF policy requires information that is unavailable in hospitals on a real-time basis.

*Key words:* queueing systems; heterogeneous servers; healthcare; hospital routing policies; fairness; quality- and efficiency-driven regime; asymptotic analysis

*History:* Received September 25, 2010; accepted October 1, 2011, by Assaf Zeevi, stochastic models and simulation. Published online in *Articles in Advance* March 9, 2012.

## 1. Introduction

Operations research methodologies, and queueing theory in particular, have generated valuable insights into operational strategies and practices, thus leading to solutions of significant problems in healthcare systems. In concert with this state of affairs, we analyze patient flow in hospitals: specifically, we focus on the emergency department (ED) and its interface with four internal wards (IWs) in a large Israeli hospital; we refer to it as Anonymous Hospital. Two operational problems could arise in this process: patients' waiting times in the ED for a transfer to the IWs could be long, and patient routing to the wards need not be fair, as far as workload allocation among the wards is concerned. In contrast to the majority of studies that address the problem of long waiting times in EDs, we explore the process of patient allocation to wards from the fairness perspective. In search of an allocation protocol that is fairness sensitive, we model the "ED-to-IW" process as a queueing system with heterogeneous server pools: the pools represent the wards, and

servers are beds. Within this modeling framework, we compare several routing strategies, thus identifying the one most appropriate for a hospital setting.

### 1.1. Motivation

Two conclusions can be drawn from Anonymous Hospital data (see Table 1 in §2.3). First, the fastest and smallest ward (ward B) is subject to the highest load: it experiences the highest number of patients per bed per month, with bed occupancy that is comparable to the other wards. The reasons behind the high turnover rate of ward B are superior managerial and staff practices, as well as (medically justified) varying policies for patients release, which results in significantly shorter average length of stay (ALOS). Importantly, a shorter ALOS does not come at the cost of inferior medical care. Indeed, the level of return rates (within three months) is comparable across wards. Hence, one could argue that the allocation process of patients from the ED to the IWs is not fair from the point of view of medical staff, a problem that has been observed in other Israeli hospitals as well.

The second empirical conclusion is that some key operational characteristics of the ED-to-IW process coincide with the features that characterize moderate-to large-scale queueing systems in the quality-and efficiency-driven (QED) regime (Erlang 1948, Halfin and Whitt 1981, Gans et al. 2003). Such a regime achieves, simultaneously, high levels of operational service quality (ED waiting times significantly shorter than service durations, specifically ALOS) and resource efficiency (high bed occupancy).

## 1.2. Contributions

Our paper focuses on modeling the ED-to-IW process, which is a key phase of patient flow in hospitals. More broadly, in 2004, the Joint Commission on the Accreditation of Healthcare Organizations (JCAHO 2004) set a standard (LD.3.10.10) for patient flow leadership. The standard requires that healthcare leaders “develop and implement plans to identify and mitigate impediments to efficient patient flow throughout the hospital.” It amplifies the need to identify the critical factors that impact patient flow, with the ultimate goal of designing and implementing policies, processes, and procedures that track, monitor, and improve patient flow throughout hospitals.

Although we concentrate on the impact of patient routing on staff fairness, our work should be viewed within the broader context of improving staff efficiency by creating an appropriate incentives structure. Operational policies that are not perceived to be fair could internalize inefficiencies, i.e., create situations where good work leads to more work. On the other hand, fair policies not only reward staff members with the best practices, but also promote the adoption of such practices. Thus, although fair policies come at a certain cost in the short run, in the long run they are likely to improve overall system efficiency.

Our contributions are as follows:

- Based on empirical data from Anonymous Hospital, we argue that an inverted- $V$  queueing model in the QED regime is appropriate for describing the ED-to-IW process, which exhibits multiscale behavior (e.g., IW lengths of stay being much longer than ED waiting times).
- We quantify operational fairness toward medical staff by means of ratios that take into account bed occupancy levels and bed *turnover rates* (the average number of admitted patients per bed per unit of time)—two main metrics as far as workload of medical staff is concerned. These metrics serve as operational proxies for fairness, which is an intricate concept. (Operational proxies relate to operations, are easy to measure, and they approximate notions that are hard to quantify.)

- Although routing algorithms that take into account fairness have been considered in the literature, we propose a practical routing algorithm (randomized most idle, or RMI) that is suitable for hospitals where only limited and partial information is available for patient routing; that is, we consider the role of information in the performance of routing policies.

- The proposed algorithm, RMI, is analyzed and compared to known algorithms. Our results demonstrate that RMI achieves a (long-run) level of fairness toward medical staff that is the same as that of algorithms that require more information to operate. Moreover, within our narrow notion of fairness, based on occupancy and turnover rates, an extension of our algorithm, weighted RMI, is flexible enough to achieve any desired fairness.

- Furthermore, our results yield important insights on differences between various routing algorithms, differences that arise at the subdiffusion (room-level) scale. In particular, these insights reveal that instantaneous imbalance of workload across the IWs is limited to just a few rooms of patients. Moreover, in the course of just a few days this imbalance (if present) disappears due to time averaging; i.e., empirical workload time averages converge to the desired long-run averages.

## 1.3. Brief Literature Review

There exists a vast amount of research on health-care systems in numerous scientific fields including operations research. We mention here the most relevant to the present work; additional related references can be found in Tseytlin (2009), Bekker and de Bruin (2010), de Bruin et al. (2010), Kc and Terwiesch (2009), González and Herrero (2004), de Véricourt and Jennings (2011), and McManus et al. (2004), as well as in each of the papers mentioned below. Readers are also referred to Green (2004, 2008), who describes the general background and issues involved in hospital capacity planning.

Research papers (e.g., Litvak et al. 2001) and popular articles (e.g., *New York Times* 2002) both recognize the importance of ED proper functioning and the consequences of its overcrowding. Patient flow from the ED to other medical units in the hospital (not just the IWs) has received special attention (even becoming the subject of a novel; Wright and King 2006). In fact, it has been acknowledged as a major trigger of ambulance diversion (Litvak et al. 2001), which is of great concern. Ramakrishnan et al. (2005) construct a two-time-scale model for a hospital system, where the wards operate on a time scale of days and are modeled by a discrete time Markov chain, and the ED operates on a much faster time scale and is modeled by a continuous time Markov chain. With the help

of this model, Ramakrishnan et al. (2005) estimate the expected occupancy of the wards and the probability of each ward reaching its capacity. The setup in Ramakrishnan et al. (2005) corresponds to ours in that the ED-to-IW process also operates in several scales. These scales arise naturally from our asymptotic analysis (QED regime). Specifically, length of stay in wards, which correspond to our service times, are naturally measured in days, whereas waiting in the ED for transfer to the wards is naturally measured in hours, which corresponds to our waiting times.

Relevant analyses of queueing systems with heterogeneous servers date back to the slow server problem (Rubinovitch 1985a). Initially, the focus was on finding the best operating policy to minimize the steady-state mean sojourn time of the customers in the system, which is equivalent to minimizing the long-run average number of customers in the system, due to Little's law. Under such criteria, it is preferable to use the faster servers more than the slower servers. In fact, under some circumstances, it is even advantageous to remove the very slow servers and thus reduce sojourn time. This slow server phenomenon is addressed, for the case of two heterogeneous servers and a single queue, by Larsen and Agrawala (1983), Lin and Kumar (1984), Rubinovitch (1985a, b), and Stockbridge (1991); and the general multiservers case is addressed by Cabral (2005) (under the random assignment (RA) policy, which routes a customer to one of the idle servers at random). Later, Cabral proved in Cabral (2007) that, for any two servers, the fraction of time that the faster among them is busy is smaller than that of the slower one, and the effective service rate of the faster server is higher than that of the slower one. Except for Cabral (2007), none of the studies mentioned above touches on the issue of *fairness* toward servers.

In the context of large-scale systems, Armony (2005) analyzed the fastest-servers-first (FSF) routing policy that assigns customers to the fastest available pool. She shows that FSF is asymptotically optimal (within the set of all nonpreemptive, nonanticipating first-come, first-served (FCFS) policies), in the sense that it stochastically minimizes the stationary queue length and waiting time as the arrival rate and number of servers grow large; Armony and Mandelbaum (2012) extended this result to accommodate abandonments. Yet, under the FSF policy, asymptotically only the slowest servers have any idle time. This is obviously unfair toward the fast servers (which get “punished” for being fast by working more), and it gives them an incentive to slow down—an undesirable result for the system as a whole. Thus, there exists a trade-off between operational optimality for the system and fairness toward servers (Armony and Ward 2010).

There is ample research aimed at achieving fairness to customers; see, for example, references in Tseytlin (2009). However, to our best knowledge, only recently, Atar (2008) was the first to deal with the operational fairness-toward-servers issue. He studied a single-server pools model in the QED regime, where the number of servers and their service rates are independent and identically distributed (i.i.d.) random variables, under a policy that routes an arriving customer to the server that has been idle for the longest time among all idle servers (in a deterministic environment, this policy is called longest-idle server first (LISF) by Armony 2005). Armony and Ward (2010) extended LISF routing to longest-weighted-idle-server-first (LWISF) routing, and proposed a threshold policy that asymptotically (in the QED regime) outperforms LWISF while achieving the same target idleness ratios (IRs). Atar et al. (2011) proposed the longest-idle-pool-first policy, that routes a customer to the pool with the longest cumulative idleness among the available pools; in the QED regime, this policy is shown to balance cumulative idleness among the pools. Most of these papers examine transient system behavior by establishing weak convergence process-level results. In contrast, our work deals with steady-state behavior.

Fairness (or justice, or equity) is a well-researched area in the behavioral sciences (Colquitt et al. 2001). We shall mention some relevant work later in §2.2.

#### 1.4. Organization

The paper is organized as follows. In the next section we describe the process for routing patients from an ED to IWs, which reveals a fairness problem in this process. Our queueing model, as well as the QED regime, are introduced in §3. Fair routing algorithms are described in §4. Section 5 contains our theoretical results and a related discussion, including managerial implications. Concluding remarks appear in §6. Readers are referred to the online appendices (Mandelbaum et al. 2011) for a description of the routing algorithm currently implemented at Anonymous Hospital, the state of affairs in five other hospitals, technical proofs, a table of notation, and a list of acronyms.

## 2. Patient Routing

We study patient flow from the ED to the IWs in hospitals. Our research site is Anonymous Hospital, which is a large Israeli hospital with approximately 1,000 beds, 45 medical units, and approximately 75,000 patients hospitalized yearly. Among its variety of units, it has a large ED with an average arrival rate of 200–300 patients daily, who occupy up to 50 beds, and five IWs, which we denote from A to E. The ED is divided into two major subunits: internal and trauma

(the latter being surgical and orthopedic patients). An internal patient, whom the ED decides to hospitalize, is directed to one of the five IWs according to a certain routing policy—this routing process is the focus of our research.

Departments of internal medicine are responsible for catering to a wide range of internal disorders, providing inpatient medical care to thousands of patients each year. Wards A–D are more or less similar in their medical capabilities—each can treat multiple types of patients. Ward E, on the other hand, treats only “walking” patients, and routing to it differs from routing to the other wards. In our study we thus concentrate on the routing process to wards A–D only. The existence of multiple wards with similar medical capabilities is common in Israeli hospitals and can be attributed to various factors, for example, (i) there exist constraints in terms of physical space, e.g., wards can be located on different floors of a building; (ii) the existence of multiple wards implies the existence of multiple positions of ward managers (these positions are used by the hospital management to attract top-performing doctors); and (iii) informal research in Anonymous Hospital suggests that health-care operations could exhibit diseconomies of scale (this would, of course, discourage the formation of a single superward).

### 2.1. The Routing Process

The decision of routing a to-be-hospitalized patient to one of wards A–D is supported by a computer program, referred to as the “Justice Table” at the hospital. As its name suggests, the algorithm’s goal is to make the patient allocation to the wards *fair*, by balancing the load among the wards. Prior to routing, patients are classified into three categories, according to the complexity of treatment: ventilated, special care, and regular. The program accepts a patient’s category as an input parameter and returns a ward for the patient (A–D) as an output. For each patient category, there is round-robin order among the wards, while accounting for the size of each ward by allocating fewer patients to a smaller ward: for example, two patients to B per three patients to A, reflecting a 30:45 bed ratio. The Justice Table does not take into account the *actual* number of occupied beds and the patient discharge rate. In Mandelbaum et al. (2011), we provide a brief history of the Justice Table and a more detailed description of the ED-to-IW process.

We recognize a problem in the process of patient routing from the ID to the IWs: patient allocation to the wards does not appear to be fair. In what follows, we first examine the notion of fairness and then discuss the specifics present at Anonymous Hospital.

REMARK 1. In addition to analyzing the ED-to-IW process in one specific hospital, we examined how

this process is managed in five other Israeli hospitals (see Mandelbaum et al. 2011). In particular, we learned about the routing policies that are being used, and how successful they are in terms of delays and fair allocation. To this end, we used a questionnaire that included both qualitative (detailed description of the process, fairness considerations) and quantitative (operational measures of the ED and IWs: capacity, ALOS, waiting times) questions. Our study revealed that unbalanced loads on the wards due to heterogeneity of ALOS are common in all our surveyed hospitals.

### 2.2. Fairness

One can analyze fairness toward *patients* (customers) and fairness toward *wards*—the latter covering medical and nursing staff (servers). There is ample literature on measuring fairness in queues from the customer point of view (e.g., see Avi-Itzhak and Levy 2004, Larson 1987, Rafaeli et al. 2002). Various aspects are investigated (for example, single queue versus multiple queues, or FCFS versus other queueing disciplines), but all agree that the FCFS policy is typically essential for justice perception. Consequently, customer satisfaction in a single queue is higher than in multiple queues (Larson 1987), and waiting in a multiqueue system produces a sense of lack of justice even when no objective discrimination exists (Rafaeli et al. 2002). The situation could be different in invisible queues (e.g., call centers) and healthcare queues (e.g., EDs). In the latter, clinical priority naturally dominates FCFS justice.

REMARK 2. We note that patient “service” in a ward actually starts prior to the physical arrival of the patient to the ward. Indeed, we observed that the ward, once informed about a to-be-admitted patient, starts preparing for this *specific* patient: different patients, even if they fall under the same classification, might require different preparations. This sometimes leads to the following scenario. Suppose that a decision for hospitalization of patient X was made prior to a decision of hospitalization of patient Y (assuming both of them are clinically similar); say patient X is directed to ward A and patient Y to ward B. Now, suppose that ward B becomes ready to physically admit the patient earlier than ward A—hence patient Y joins a ward *before* patient X, although Y “arrived” later (Elkin and Rozenberg 2007). In addition to this need for a ward to prepare for the patient, the hospital staff refrains from modifying patient-ward assignments also for a psychological reason: a patient awaiting hospitalization (as well as accompanying individuals) experiences high levels of stress as is—one thus does not wish to aggravate stress by changing original ward assignments (Elkin and Rozenberg 2007).

In our work, we focus on the process of patient assignment to wards, as opposed to the physical process of patient transfer from the ED to the IWs. We refer to the former as the ED-to-IW process. (In this process, deviations from FCFS do not raise fairness problems among patients since the assignment queue is typically invisible to them.)

The literature on justice from the server point of view is concerned with *Equity Theory* (Huseman et al. 1987), according to which workers perceive their justice by comparing ratios of outputs from the job to inputs to the job. Specifically, if the output/input ratio of an individual is perceived to be unequal to others, then inequity exists. The larger the inequity the individuals perceive, the more uncomfortable they feel and the harder they work to restore equity (Huseman et al. 1987). Ben-Zrihen et al. (2007) showed that in customer service centers, servers' equity perception has a positive influence on their performance and job satisfaction. References to additional studies on the importance of perceived justice among employees can be found in Armony and Ward (2010).

Our anchor point is the survey reported in Elkin and Rozenberg (2007), in which the staff (nurses, doctors, and administration) were asked to grade the extent of fairness in different routing policies. When discussing fairness with ward staff, the consensus was that each nurse/doctor should have the same workload as others. Seemingly, this is the same as saying that each nurse/doctor should take care, at any given time, over an equal number of patients (assuming homogeneous customers, for simplicity). Because the number of nurses and doctors is usually proportional to standard capacity, this criterion is equivalent to keeping *occupancy levels* of beds equal among the wards.

Note that, by Little's law,  $\rho = \gamma \times \text{ALOS}$ , where  $\rho$  is the average occupancy level, and  $\gamma$  is the bed turnover rate. Thus, if one maintains ward occupancies equal then wards with shorter ALOS will have a higher turnover rate—admit more patients per bed—which gives rise to additional fairness concerns. Indeed, the load on staff is not spread uniformly over a patient's stay, because treatment during the first days of hospitalization requires much more time and effort from the staff than in the following days (Elkin and Rozenberg 2007); moreover, patient admissions and discharges significantly consume doctors' and nurses' time and effort as well. Thus, even if occupancy among wards is kept equal, the ward admitting more patients per bed ends up having a higher load on its staff. For these reasons, a natural alternative fairness criterion is balancing the turnover rate, or the *flux*—namely, the number of admitted patients per bed per unit of time (for example, per month)—among the wards.

One can also combine utilization and flux to produce a single workload measure; fair routing would maintain this measure equal across wards. For example, load, experienced by medical staff in a ward, can be roughly divided into two parts: load associated with treating hospitalized patients (quantified by utilization) and load due to patient admissions/discharges (quantified by flux). For example, a single (objective) workload measure for a nurse, based on a linear combination of utilization and flux, was proposed by Tseytlin (2009, §7.2):

$$\frac{\gamma_i N_i T_i^\gamma + \rho_i N_i T_i^\rho}{n_i}, \quad (1)$$

where  $\gamma_i$  is the flux,  $\rho_i$  is the occupancy level,  $N_i$  is the number of beds in the ward,  $T_i^\gamma$  is the average amount of time required from a nurse to complete one admission plus discharge,  $T_i^\rho$  is the average time of treatment required by a hospitalized patient per unit of time, and  $n_i$  is the number of nurses in the ward; all quantities refer to a given ward  $i$ .

### 2.3. The Setting of Anonymous Hospital

Although wards A–D in Anonymous Hospital provide similar medical services, they do differ in their operational characteristics (see Table 1), which we now elaborate on. First of all, each medical unit is characterized by its *capacity*. A ward's capacity is measured by its number of beds (standard static capacity) and number of service providers—doctors, nurses, administrative staff, and support staff (dynamic capacity). The “maximal” static capacity stands for the standard static capacity plus extra beds, which can be placed in corridors during overloaded periods. It is convenient to introduce notions of a subward and a patient room; these will play a role in a discussion of our results. Wards consist of subwards, which, in turn, are made up of physically collocated rooms. There a few (two to three) subwards in

**Table 1** Internal Wards Operational Measures

	Ward A	Ward B	Ward C	Ward D
ALOS <sup>a</sup> (days)	6.5 (±0.19)	<b>4.5 (±0.15)</b>	5.4 (±0.22)	5.7 (±0.18)
Mean occupancy level (%)	97.8	94.4	86.8	91.1
Mean # patients per month	205.5	187.6	210.0	209.6
Standard (maximal) capacity (# beds)	45 (52)	30 (35)	44 (46)	42 (44)
Mean # patients per bed per month	4.58	<b>6.38</b>	4.89	4.86
Return rate (within 3 months) (%)	16.4	17.4	19.2	17.6

*Notes.* Data refer to the period May 1, 2006–October 30, 2008 (excluding January–March 2007, when ward B was in charge of an additional subward). The data cover 16,947 admissions in total.

<sup>a</sup>The level of confidence for average length of stay (ALOS) is 95%.

**Table 2** Numbers of Admitted Patients to Wards A–D for Each Patient Category

IW\Patient type	Regular	Special care	Ventilated	Total
Ward A	2,316 (50.3%)	2,206 (47.9%)	83 (1.8%)	4,605 (25.2%)
Ward B	1,676 (43.0%)	2,135 (54.7%)	90 (2.3%)	3,901 (21.4%)
Ward C	2,310 (49.9%)	2,232 (48.2%)	88 (1.9%)	4,630 (25.4%)
Ward D	2,737 (53.5%)	2,291 (44.8%)	89 (1.7%)	5,117 (28.0%)
Total	9,039 (49.5%)	8,864 (48.6%)	350 (1.9%)	18,253

*Note.* Data refer to the period May 1, 2006–September 1, 2008 (excluding the months January–March 2007, when ward B was in charge of an additional subward).

a ward, and a subward consists of several (three to four) rooms.

In our hospital IWs, the dynamic capacity during a particular shift is determined proportionally to the static capacity (see, however, discussions on the appropriateness of such “proportional” staffing in de Véricourt and Jennings 2011, Green 2008, Yom-Tov 2007). In particular, an IW beds-to-nurses ratio in morning shifts is 5:1 or 6:1 (depending on the number of “intensive care” beds in a ward); during night shifts, the ratio is 8:1 or 9:1. Note that the number of nurses is determined by the number of beds in a ward rather than the number of occupied beds (patients). Hence a unit’s operational capacity can be characterized by its number of beds only—denoted as its *standard* capacity.

#### 2.4. Fairness at Anonymous Hospital

Medical units are further characterized by various performance measures: *operational* (average bed occupancy level, ALOS, waiting times for various resources, number of patients admitted or released per bed per time unit (*flux*)) and *quality* (patients’ return rate, patients’ satisfaction, mortality rate, etc.). Note that occupancy levels and flux are calculated relatively to wards’ standard capacities. (Thus, occupancy can exceed 100%.) Comparing the two basic measures, ward capacity and ALOS, we observe in Table 1 that the wards differ in both. Indeed, ward B is significantly the smallest and the “fastest” (shortest ALOS) among wards A–D. We observe that the mean occupancy rate in this ward is high (comparable to that in wards A and D; higher than in ward C). In addition, the number of patients hospitalized per month in this ward is about 90% of those hospitalized per month in the other wards, although its size is merely about two-thirds of the others. As a consequence, the flux in ward B (6.38 patients per bed per month) is significantly the highest. Because nurses and doctors in Anonymous Hospital are assigned to particular wards and are salaried workers (as opposed to hourly workers), the load on the ward B staff is hence the highest. Because the staff-to-beds ratio is fixed, if ALOS is kept constant, the occupancy

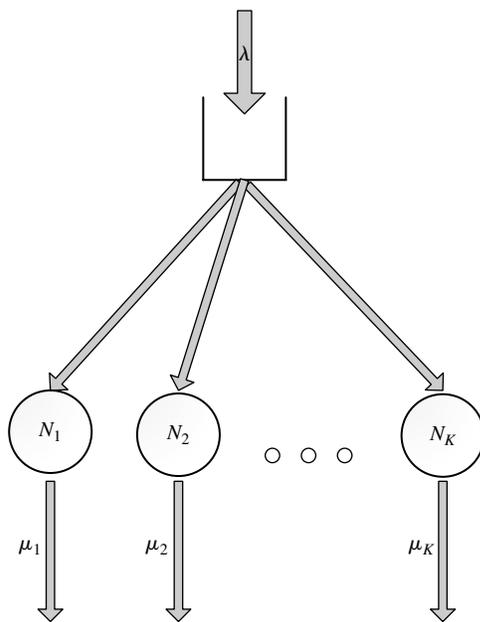
rate in a ward serves a measure of the load on the staff in that ward.

Short ALOS could be caused by multiple reasons. For example, it can result from a superior efficient clinical treatment or a liberal (versus conservative) release policy; a clinically too-early discharge of patients is clearly undesirable. One possible (and accessible) quality measure of clinical care is patients’ *return rate* (proportion of patients who are rehospitalized in the IWs within a certain period of time—in our case, three months). In Table 1 we observe that the return rate in ward B does not differ significantly from the other wards. Also, patients’ satisfaction level in surveys—another measure of care quality—does not differ in ward B (Elkin and Rozenberg 2007). Furthermore, the difference in ALOS across wards could be due to different wards treating different types of patients and/or performing different types of procedures. However, our data do not support this hypothesis. In fact, according to Table 2, ward B handles a disproportionate share of special-case and ventilated patients, patients that require longer ALOS on average. Rather, a shorter ALOS in ward B can be attributed to superior staff practices. We conclude that the most efficient ward, instead of being rewarded, is exposed to the highest load; hence, patient allocation appears unfair, as far as the wards are concerned. Increasing fairness in the routing process is expected to increase staff satisfaction, as well as provide incentives for improved care and cooperation.

### 3. Model Formulation

We model the ED-to-IW process as a queueing system with heterogeneous pools of i.i.d. servers. Arrivals to the system are patients to be hospitalized in the IWs; pools represent the IWs, which indeed have different service rates ( $1/\text{ALOS}$ ); and the number of servers in each pool corresponds to the number of beds in each ward. To create a tractable system, we assume that arrivals to the wards occur according to a Poisson process, and LOSs in wards are exponentially distributed. (Both assumptions are important for analytical tractability, but they are inaccurate reality-wise; see Armony et al. (2011) for empirical findings on

Figure 1 Inverted-V System



arrivals and LOS). Although, as remarked in §2.1, patients to be hospitalized in the IWs are classified into several categories, we analyze here a single customer class model. This, as well as our distributional assumptions, certainly could present a limitation for application of our theoretical results, but we are still able to draw useful insights about fair routing in hospitals.

### 3.1. The Inverted-V Model

Consider the queueing system shown in Figure 1. This  $\wedge$ -model, or inverted-V model (in the terminology of Armony 2005), consists of  $K$  server pools: pool  $i$  has  $N_i$  i.i.d. servers, each with exponential service times of rate  $\mu_i$  (namely, service rates are equal within each pool but vary among the pools). The total number of servers in the system is  $N = \sum_{i=1}^K N_i$ . Upon arrival, a customer is routed to one of the available pools (if it has one or more idle servers) or joins a centralized queue of infinite capacity if all the servers at all pools are busy. Homogeneous customers arrive according to a Poisson process with rate  $\lambda > 0$ . Each customer requires a service that can be provided by any of the servers, and each server can serve only one customer at a time. The queueing discipline is FCFS, nonpreemptive (the service of a customer can not be interrupted once started), and work conserving (there are no idle servers whenever there are customers awaiting service in the queue). In addition, we assume that all interarrival and service times are statistically independent.

REMARK 3. Our model is centered around beds rather than personnel. In the setup of Anonymous

Hospital (see §§2.3 and 2.4), there exists a direct connection between the number of patients (occupied beds) and staff workload in a ward, a key feature that is utilized in our model. It is possible to conceive of an alternative, more complex model that focuses on the nursing staff directly; nevertheless, as explained in §6.2, our simpler model still captures the essentials of fair routing, if looked at through the appropriate “lenses.”

### 3.2. The QED Asymptotic Regime

In this section, we formally introduce the QED regime in the context of the inverted-V model. We provide some justification for its applicability in modeling the ED-to-IW process as well.

The QED regime was first discovered by Erlang (1948); it was mathematically formalized by Halfin and Whitt (1981), hence it is often referred to as the Halfin-Whitt regime. The regime can be informally characterized in terms of any one of the following conditions: a (steady-state) system with a large volume of arrivals (demand) and many servers (supply or capacity) is operating in the QED regime if (i) the delay probability is neither near 1 nor near 0, (ii) its waiting time is one order of magnitude shorter than service time (e.g., seconds versus minutes in call centers, hours versus days in our case), or (iii) its total service capacity is equal to its demand up to a safety capacity, which is of the same order of magnitude as the square root of the demand. Characterizations (i) and (ii) relate to the quality aspect, and characterization (iii) points at high server efficiencies; thus, the QED regime achieves high levels of both service quality and system efficiency by carefully balancing between the two (Armony 2005, Gans et al. 2003). The suitability of the QED regime to the ED-to-IW process was studied in detail in Armony et al. (2011). Here we emphasize the following relevant empirical facts:

- The number of servers (beds) in each pool (ward) is approximately 30–50 (Table 1)—the system is large enough for QED approximations to apply (Borst et al. 2004, Zhang et al. 2012, Yom-Tov 2010).
- Servers’ utilization (bed occupancy) is above 85% (Table 1).
- Waiting times are indeed an order of magnitude shorter than service times: hours versus days. Observe that a waiting time in our inverted-V model corresponds to a patient waiting in the ED due to lack of available beds in IWs only, i.e., bed allocation time—the amount of time from the decision to hospitalize the patient until a bed becomes available in one of the IWs. In particular, if a bed is available at the time of the hospitalization decision, the corresponding bed allocation time is zero. However, estimating bed allocation times is difficult because they are not tracked by Anonymous Hospital’s information systems. Hence, we estimated the time between

the hospitalization decision and the time of the first procedure in the IW (see Mandelbaum et al. 2011). This time serves as an upper bound on bed allocation time. Because the latter is a major component of the total waiting time in the ED, we conclude that both quantities are in the order of hours. A service time in the inverted- $V$  model corresponds to the interval from the time when a patient occupies a bed until the time the patient releases the bed (or, more precisely, until the bed becomes available next).

- The probability of encountering an available bed in the designated ward, upon hospitalization decision (relatively to the wards' standard capacities), is estimated to be 43%, 48%, 76%, and 55%, for wards A–D, respectively. The probability of encountering an available bed in any of the wards is approximately 84%; this estimate is based on a long period of time, which includes lightly loaded periods (the arrival process is nonstationary). If one considers only highly loaded winter months (November to April), the probabilities of encountering an available bed in wards A–D are 29%, 35%, 64%, and 45%, respectively; the probability of encountering an available bed in any of the wards is 75% (this probability is the lowest in January—59%).

**REMARK 4.** The probability of being admitted to a ward immediately (or within a short time) after the hospitalization decision is much smaller than the probability of encountering an available bed (see Figure 7 in Mandelbaum et al. 2011). Indeed, during the period May 1, 2006, to October 30, 2008 (excluding the months January–March of 2007), only 2.7% of the patients were admitted to an IW within 15 minutes from their assignment to a ward. This fact is consistent with the *efficiency-driven* regime. We further note that, in evening shifts (when most patients are admitted; Armony et al. 2011), there usually are one or two doctors in each ward; hence, one expects that they operate under high loads. In that case, the probability of encountering an available doctor, which is a prerequisite for being admitted to a ward, should indeed be at ED levels (but we do not have any data on staff availability). We thus have two parallel queueing systems: beds, which are QED, and medical staff, who are ED. As already indicated, we focus on the former.

We use the following scaling, suitable for the inverted- $V$  model. Consider a sequence of systems, indexed by  $\lambda$  (to appear as a superscript), with increasing arrival rates  $\lambda \rightarrow \infty$ , and increasing total number of servers  $N^\lambda$ , but with fixed service rates  $(\mu_1, \dots, \mu_K)$ . Then, the service capacity of pool  $i$  is  $c_i^\lambda = N_i^\lambda \mu_i$ , the total service capacity is  $c^\lambda = \sum_{i=1}^K c_i^\lambda$ , and the total traffic intensity is  $\rho^\lambda = \lambda/c^\lambda$ . Both  $\lambda$  and  $N^\lambda$  tend to infinity simultaneously, in a way that two limiting relations are satisfied. First, for large  $\lambda$ ,

each pool has a nonnegligible fraction of the total capacity:

$$c_i^\lambda/c^\lambda \rightarrow a_i, \quad \text{as } \lambda \rightarrow \infty, \quad (\text{C1})$$

where  $a_i > 0$  ( $i = 1, 2, \dots, K$ ), and  $\sum_{i=1}^K a_i = 1$ . The scalar  $a_i$  is the limiting proportion of the service capacity of pool  $i$  ( $i = 1, 2, \dots, K$ ) out of the total capacity. Second, it is convenient to define a scaling parameter,

$$\nu^\lambda := \lambda/\hat{\mu},$$

where  $\hat{\mu}$  is the arithmetic-mean service rate,

$$\hat{\mu} := \sum_{i=1}^K a_i \mu_i;$$

it is appropriate to think of  $\nu^\lambda$  as an effective system size. Then, the following condition is assumed to hold:

$$\sqrt{\nu^\lambda} (1 - \rho^\lambda) \rightarrow \delta > 0, \quad \text{as } \lambda \rightarrow \infty. \quad (\text{C2})$$

This limit implies  $\rho^\lambda \rightarrow 1$  (high utilization), as  $\lambda \rightarrow \infty$ , and is (asymptotically) equivalent to the classical square-root safety staffing rule: the total service capacity,  $c^\lambda = \lambda + \delta\sqrt{\lambda\hat{\mu}} + o(\sqrt{\lambda})$ , is equal to the arrival rate,  $\lambda$ , plus a square-root safety capacity,  $\delta\sqrt{\lambda\hat{\mu}}$ , where  $\delta$  is some quality-of-service parameter (the larger the value of  $\delta$ , the higher the service quality); here,  $\delta$  is a unitless quantity, whereas  $c^\lambda$ ,  $\lambda$ , and  $\sqrt{\lambda\hat{\mu}}$  are measured in the same units (e.g., patients/week). Note that fluctuations of the arrival process from its mean ( $\lambda$ ) are proportional to  $\sqrt{\lambda}$ . With  $\mu$  being the harmonic-mean service rate,

$$\mu^{-1} := \sum_{i=1}^K a_i \mu_i^{-1},$$

(C1) and (C2) then imply

$$\frac{\lambda}{N^\lambda} = \frac{\lambda}{\sum_{i=1}^K c_i^\lambda/\mu_i} \rightarrow \mu,$$

as  $\lambda \rightarrow \infty$ . In view of this, (C2) can be rewritten as

$$\sqrt{N^\lambda} (1 - \rho^\lambda) \rightarrow \beta := \delta\sqrt{\hat{\mu}/\mu},$$

as  $\lambda \rightarrow \infty$ . Finally, we define  $q_i$  to be the limiting fraction of pool  $i$  servers, out of the total number of servers:

$$\frac{N_i^\lambda}{N^\lambda} \rightarrow \frac{a_i}{\mu_i} \mu := q_i, \quad (2)$$

as  $\lambda \rightarrow \infty$ ; the limit is due to (C1) and (C2). Because  $a_i > 0$ , for all  $i$ , we also have that  $q_i > 0$ , for all  $i$ , i.e., the pools are of comparable sizes. Clearly,  $\sum_{i=1}^K q_i = 1$ , and  $\sum_{i=1}^K q_i \mu_i = \mu$ . Therefore, one can interpret  $\mu$  as the (limiting) average service rate of a server in the system, whereas  $\hat{\mu}$  is the (limiting) average service rate at which customers receive service. The two quantities differ because faster servers serve more customers. (Note that  $\mu \leq \hat{\mu}$ , with equality if and only if all  $\mu_i$ 's are equal to each other, in which case  $\beta = \delta$  and  $\nu^\lambda = N^\lambda$ .)

## 4. Fair Routing

Before introducing formally two criteria of fairness within the context of the inverted- $V$  model, we need some notation. Denote by  $I_i^\lambda$  the long-run (steady-state) number of idle servers at pool  $i$  ( $i = 1, 2, \dots, K$ ); each  $I_i^\lambda$  is a random variable that attains values in  $\{0, 1, \dots, N_i^\lambda\}$  ( $i = 1, 2, \dots, K$ ). Let  $\rho_i^\lambda := 1 - \mathbb{E}I_i^\lambda / N_i^\lambda$  be the mean long-run (steady-state) occupancy rate in pool  $i$ . As the servers within each pool are symmetric,  $\rho_i^\lambda$  also stands for the utilization of servers in pool  $i$ —the fraction of time that each server is busy in the long run (steady state). Finally, we denote by  $\gamma_i^\lambda$  the average flux through pool  $i$  (average number of service completions per pool  $i$  server per time unit):  $\gamma_i^\lambda = \mu_i \rho_i^\lambda$ , by Little’s law. Clearly,  $\gamma_i$  stands also for the average effective service rate of a server in pool  $i$ .

When analyzing fairness toward servers, we consider the following two criteria: *occupancy* balancing and *flux* balancing. The server’s utilization, or, equivalently, the pool’s occupancy rate, is one of the prevalent measures of the server’s workload. As the occupancy rates at all pools tend to one in the QED regime (see (C2)); we thus compare the ratio between the proportions of *idle* servers in the pools  $(1 - \rho_i^\lambda) / (1 - \rho_j^\lambda)$ , referred to as the *idleness ratios*. The closer these ratios are to unity, the more balanced the routing is, according to the occupancy-balancing criterion. The second criterion takes into account the additional measure of workload, the average “flux” through the pools, namely, the number of customers served by a server per time unit. Hence, our second criterion is based on the *flux ratios*  $\gamma_i / \gamma_j$ . The closer this ratio is to unity, the more balanced the routing is, according to the flux-balancing criterion. Note that, in the QED regime,  $\gamma_i / \gamma_j \rightarrow \mu_i / \mu_j$  as  $\lambda \rightarrow \infty$  for any work-conserving routing algorithm; that is, pools with higher service rates experience higher flux. Hence, based on the discussion in §2.2, it is appropriate to have lower utilization for pools with higher service rates. Equivalently, if the flux ratio for two pools is greater than one (because of the difference in the service rates), then the idleness ratio should be greater than one as well. The algorithms we discuss in the next subsection all achieve this goal.

### 4.1. Routing Algorithms

In this subsection, we describe three routing algorithms. All three are work conserving, a choice that is dictated by our goal to reduce queue length (or, equivalently, waiting time by Little’s law) rather than the number-in-system (sojourn time). In the ED-to-IW setting, the objective is to minimize the “queue” for transfer to the wards, thus reducing the overload on the ED. Work conservation is discussed further in the remark at the end of the present subsection.

Next, we describe the three routing algorithms and their implications on server fairness. Let  $I_i^\lambda(t) \in \{0, \dots, N_i^\lambda\}$  denote the number of idle servers in pool  $i$  ( $i = 1, 2, \dots, K$ ) at time  $t$ . When there are no customers awaiting service at time  $t$ , the vector  $(I_1^\lambda(t), \dots, I_K^\lambda(t))$  specifies the state of the system. However, when the waiting queue is not empty, an additional variable is needed to specify the number of customers awaiting service; let  $Q^\lambda(t) \in \{0, 1, \dots\}$  be the number of customers awaiting service at time  $t$ . It is convenient to define  $I^\lambda(t) \in \{\dots, N^\lambda - 1, N^\lambda\}$  as (jointly) the total number of idle servers awaiting customers/number of customers awaiting servers, at time  $t$ :

$$I^\lambda(t) = \sum_{i=1}^K I_i^\lambda(t) - Q^\lambda(t). \quad (3)$$

Note that  $I^\lambda(t)$  can take negative values:  $\{I^\lambda(t) = -i\}$ ,  $i \geq 1$ , indicates that there are  $i$  customers awaiting service at time  $t$ . Because of work conservation, one has  $\sum_{i=1}^K I_i^\lambda(t) = (I^\lambda(t))^+$  and  $Q^\lambda(t) = (I^\lambda(t))^-$ , where  $x^+$  and  $x^-$  denote the positive and negative parts of  $x$ , respectively.

**4.1.1. LISF Routing.** The LISF policy routes a customer to a server that has been idle for the longest time among all idle servers. This policy is commonly used in call centers and considered to be fair (Armony and Ward 2010). It was first analyzed (in the QED regime) by Atar (2008) in the context of a single-pool system in a random environment (service rates were taken to be i.i.d. random variables). Armony and Ward (2010) analyzed LISF routing in the inverted- $V$  model with two ( $K = 2$ ) pools. Informally, they show that, for large  $\lambda > 0$ , the LISF policy maintains fixed ratios between the number of idle servers in different pools, whenever idle servers exists; that is, if  $I^\lambda(t) > 0$  at time  $t$ , then

$$\frac{I_i^\lambda(t)}{I_j^\lambda(t)} \approx \frac{a_i (I^\lambda(t))^+}{a_j (I^\lambda(t))^+} = \frac{a_i}{a_j}.$$

Hence, up to some technical conditions (see the discussion prior to Corollary 4.2 in Gurvich and Whitt 2009), under LISF routing one expects

$$\frac{1 - \rho_i^\lambda}{1 - \rho_j^\lambda} = \frac{\mathbb{E}I_i^\lambda}{N_i^\lambda} \Big/ \frac{\mathbb{E}I_j^\lambda}{N_j^\lambda} \rightarrow \frac{a_i q_j}{a_j q_i} = \frac{\mu_i}{\mu_j},$$

and  $\gamma_i^\lambda / \gamma_j^\lambda \rightarrow \mu_i / \mu_j$ , as  $\lambda \rightarrow \infty$ ; i.e., both the idleness and flux ratios tend to the ratios of server rates. The algorithm leads to a desirable outcome: fast servers work less (have lower utilization) but “produce” more (have higher flux).

Note that even though LISF is a “blind” policy (a policy that requires, at the time of routing decision, none or minimal information on the parameters

of the system or the system state; Atar et al. 2011), implementing the LISF policy in the hospital setting is not straightforward; namely, one must keep track not only on the number of idle servers (beds) in each pool (ward), but also the relative ordering of idle servers in terms of their idle times. The latter information is not currently available in hospitals. This fact motivates us to consider alternative routing policies that achieve the same (asymptotic) fairness toward servers but utilize less information for customer routing.

**4.1.2. IR Routing.** The IR routing policy is a way to achieve the same idleness ratio as the LISF policy, but without the information on idleness times. This policy is a special case of queues-and-idleness-ratio (QIR) policies, which were proposed and analyzed by Gurvich and Whitt (2009). They consider a generalization of the inverted- $V$  model, a parallel-server system—a service system with multiple server pools and multiple customer classes. The problem of dynamic control of such systems is often referred to as “skill-based routing” (borrowing the terminology from the world of call centers). Adopting the QIR policy to the inverted- $V$  model, IR routes a customer to the server pool with the highest idleness imbalance. The basic idea is to route each customer in such a way that the vector  $(I_1^\lambda(t), \dots, I_K^\lambda(t))$  is as close to  $(w_1(I^\lambda(t))^+, \dots, w_K(I^\lambda(t))^+)$  as possible, where the set of positive constants  $\{w_i\}$  is a priori fixed, with  $\sum_{i=1}^K w_i = 1$ . For example, if  $K = 2$  and  $w_1 = w_2 = 1/2$ , then the IR policy routes a customer to the pool with the higher number of idle servers. Specifically, at time  $t$ , a customer is routed to the pool with the index

$$\arg \max_i \{I_i^\lambda(t-) - w_i(I^\lambda(t-))^+\}.$$

Ties are broken in an arbitrary but consistent fashion; the tie-breaking rule does not impact the results at the diffusion ( $\sqrt{\nu^\lambda}$ ) scale. Gurvich and Whitt (2009) showed that (Q)IR controls drive the idleness process to the predetermined proportions  $\{w_i\}$ , in the QED regime. Following the same argument as in the LISF case, one expects

$$\frac{1 - \rho_i^\lambda}{1 - \rho_j^\lambda} = \frac{\mathbb{E}I_i^\lambda}{N_i^\lambda} \Big/ \frac{\mathbb{E}I_j^\lambda}{N_j^\lambda} \rightarrow \frac{w_i q_j}{w_j q_i},$$

and  $\gamma_i^\lambda / \gamma_j^\lambda \rightarrow \mu_i / \mu_j$ , as  $\lambda \rightarrow \infty$ . Therefore, with the IR algorithm, one can achieve the same ratio of the number of idle server (and idleness ratios) as under the LISF algorithm by simply setting  $w_i = a_i$ , because  $(w_i q_j) / (w_j q_i) = \mu_i / \mu_j$  (see (2)). Given its weights, the IR policy is a blind policy as only the information on the number of idle servers in each pool is needed. However, determining the values of  $a_i$  is far from straightforward in the hospital setting. In particular, note that  $a_i$  represents the limiting ratio between the

pool  $i$  service capacity ( $c_i^\lambda$ ) and the total service capacity ( $c^\lambda$ ). Therefore, because one must estimate the service rates to evaluate the appropriate weights, the considered version of the policy is, effectively, not blind. Recall from §2 that the capacity (number of beds) of each ward can vary with time; e.g., during the first three months of 2007, ward B was in charge of an additional subward. These facts serve as a motivation for considering yet a third routing algorithm, based on further reduced information.

**4.1.3. RMI Routing.** We now introduce the RMI routing policy. As our results in the next section indicate, the RMI policy achieves the same ratios of idleness between server pools as the LISF and IR (with  $w_i = a_i$ ) policies on the diffusion scale, and yet it requires information on neither idleness times nor on pool capacities. Under the RMI policy, a customer is assigned to one of the available pools, with probability that equals the fraction of idle servers in that pool out of the overall number of idle servers in the system (hence, the name of the policy); that is, a customer to be routed at time  $t$  is assigned to pool  $i$  with probability  $I_i^\lambda(t-) / (I^\lambda(t-))^+$ . Tseytlin (2009) observed that the RMI policy, in the inverted- $V$  model, is equivalent to the RA policy in the single-server pool model with  $N^\lambda$  servers, where  $N_i^\lambda$  servers have rate  $\mu_i$  ( $i = 1, \dots, K$ ). To illustrate this, consider the following example. Assume that a customer arrives to the system with two available pools: pool  $i$  has two available servers, and pool  $j$  has three available servers. Thus, the customer is routed to pool  $i$  with probability  $2/5$  and to pool  $j$  with probability  $3/5$ . As in pool  $i$ , there are two available i.i.d. servers, and each one of them will serve this customer with probability  $1/5$ ; similarly, each server in pool  $j$  will serve the customer with probability  $1/5$ . As a consequence, the customer is assigned to any one of the five available servers with equal probabilities.

An analytically appealing feature of the RMI policy is that, when modeled as a Markov chain in continuous time, the system is reversible (Kelly 1979) (Tseytlin 2009 conjectured that this is the only routing policy under which the inverted- $V$  system induces a reversible Markov chain). Therefore, one is able to derive its steady-state probabilities in a straightforward manner and provide an exact analysis in steady state. (Because of their complexities, LISF and IR policies were analyzed only asymptotically.) Finally, we note that, even though RMI is a randomized policy, it can be easily implemented in a hospital setting. For example, patient ID numbers can be utilized as sources of randomness (as stated by Mandelbaum et al. 2011, at least one hospital in our survey has implemented some randomized policy, based on patient ID numbers).

REMARK 5. (WHY WORK CONSERVING?). The three considered algorithms are work conserving. This assumption is not restrictive from the point of view of both hospital and patients. The goal is to reduce the waiting times in the ED so that the number of patients awaiting transfer from the ED to the IWs is minimized. Therefore, implementing a work-conserving policy is desirable because intentional server idling only increases the number of patients in the ED. Tseytlin (2009) showed that, under the waiting time criterion and RMI routing, it is not beneficial to discard slow servers, i.e., it is not desirable to intentionally idle servers. On the other hand, when the objective is to minimize the sojourn time (waiting plus service times), cases arise when it is beneficial to discard slow servers in a system with heterogeneous servers. Within the context of the well-known *slow-server* problem, this was shown for two servers (Rubinovitch 1985a) and many servers (Cabral 2005); namely, if some servers are slow enough, roughly speaking, it could turn out preferable for a waiting customer to wait for a fast server to become available rather than start service at a slow server. Thus, as far as sojourn time is concerned (but not waiting time), a non-work-conserving policy could turn out preferable.

## 5. Theoretical Results

The following theorem characterizes RMI performance, in the nonasymptotic regime (finite arrival rate  $\lambda$ ). The theorem states that when comparing two pools, servers utilization in the faster pool is lower than that in the slower pool, but the flux in the faster is higher than in the slower. Because of the symmetry of servers within each pool, this implies that, when considering any two servers, the faster between the two will work less time than the slower one, but, at the same time, the faster server will serve more customers than the slower. The result suggests, first of all, some form of fairness: faster servers are “rewarded” by working less time. In addition, operational preferences of the system are accommodated as well: more customers are served by faster servers than by slower servers. We note that Cabral (2007) proved this result, for the single-server-pool system under the RA policy, independently.

**THEOREM 1.** *In the inverted- $V$  model under the RMI policy, for any two pools  $i$  and  $j$ , if  $\mu_i > \mu_j$ , then  $\rho_i^\lambda < \rho_j^\lambda$  and  $\gamma_i^\lambda > \gamma_j^\lambda$ .*

**PROOF.** See Mandelbaum et al. (2011).  $\square$

The theorem provides an upper and lower bound on the ratio of server utilizations (assuming  $\mu_i > \mu_j$ ):  $\mu_j/\mu_i < \rho_i^\lambda/\rho_j^\lambda < 1$ . This suggests that the difference in utilizations of any two servers is more significant the more their service rates differ: for  $\mu_j \approx \mu_i$ ,

one has  $\rho_j^\lambda \approx \rho_i^\lambda$ , but as  $\mu_j$  grows smaller than  $\mu_i$ , the server-utilization ratio decreases. The bounds on the flux ratios are as follows (assuming  $\mu_i > \mu_j$ ):  $1 < \gamma_i^\lambda/\gamma_j^\lambda < \mu_i/\mu_j$ , where the second inequality follows from  $\rho_i^\lambda/\rho_j^\lambda < 1$ . The latter upper bound is important—although the fact that faster servers serve more customers contributes to system performance, one should not forget that, in certain cases, higher flux actually implies higher workload. In particular, this is the case in our ED-to-IW process, because service admissions and releases impose workload that is plausibly proportional to flux. For the RMI policy, the servers’ flux ratio  $\gamma_i^\lambda/\gamma_j^\lambda$  is bounded by the ratio of their service rates:  $\gamma_i^\lambda/\gamma_j^\lambda = (\rho_i^\lambda \mu_i)/(\rho_j^\lambda \mu_j) < \mu_i/\mu_j$  when  $\mu_i > \mu_j$ ; thus, if server rates are comparable, a faster server’s flux can not be much higher than that of a slower one.

REMARK 6. The fact that utilization decreases the faster the server gets, provides an incentive for servers to work faster, which is positive on one side but, on the other, might harm service quality, if one starts serving customers *too* fast. For example, see the study by Gans et al. (2003), who describe telephone agents that intentionally hang up on customers to maintain low average service times. Another issue is that if the slow server is not responsible for being slow (for example, a new server versus a veteran), “punishing” the server for being slow appears quite unfair. However, as noted earlier, higher flux may be considered in certain cases a “punishment” as well. Hence, to decide which servers perceive themselves as better off (the faster or the slower) or, alternatively, what servers’ incentives are (to increase or decrease his/her service rate), one must account for the servers utility functions (the ones they strive to optimize), which combine both criteria (utilization and flux); for an example of such a utility function, recall (1).

The next theorem characterizes the inverted- $V$  model under RMI routing in the QED regime. It can serve as a means for evaluating performance measures (probability of wait, expected waiting time, etc.) of the inverted- $V$  model in the QED regime. A summary of performance measures, both for finite  $\lambda$  and in the limit, as  $\lambda \rightarrow \infty$ , can be found in Mandelbaum et al. (2011). Recalling (3), let  $I^\lambda$  be the *stationary* total number of idle servers/customers awaiting service in the system with arrival rate  $\lambda$ ; the random variable  $I^\lambda$  takes values in  $\{\dots, N^\lambda - 1, N^\lambda\}$ ; by definition,  $(I^\lambda)^+ = \sum_{i=1}^{N^\lambda} I_i^\lambda$ ; and  $\{I^\lambda = i\}$  for negative  $i$  indicates that there are  $|i|$  customers awaiting service. As the system size increases ( $\lambda \rightarrow \infty$ ), the variability of demand increases as well, implying that the magnitude of  $I^\lambda$  gets larger and larger. Hence, to gain understanding of the behavior of  $I^\lambda$ , we consider its scaled version:

$$\hat{I}^\lambda := I^\lambda / \sqrt{\nu^\lambda};$$

such a scaling is typical for the QED regime and is referred to as a diffusion scaling. Informally, the theorem states that, in stationarity, there exists a dimensionality reduction in the sense that the stationary number of idle servers in pool  $i$  satisfies  $I_i^\lambda \approx a_i(I^\lambda)^+$ , for large  $\lambda$ ; i.e., idle servers are distributed across the pools proportionally, according to the relative pool capacities. Moreover, the theorem turns out to provide an explicit estimate of how close  $I_i^\lambda$  is to  $a_i(I^\lambda)^+$ . In particular, fluctuations of  $I_i^\lambda$  around  $a_i(I^\lambda)^+$  grow with  $\lambda$  at a rate  $\sqrt[4]{\nu^\lambda}$ , and thus we consider

$$\hat{I}_i^\lambda := \frac{1}{\sqrt{I^\lambda}} \left( I_i^\lambda - \frac{c_i^\lambda}{c^\lambda} I^\lambda \right), \quad i = 1, \dots, K.$$

We refer to the  $\sqrt[4]{\nu^\lambda}$  scale as a subdiffusion scale because  $\sqrt[4]{\nu^\lambda} \ll \sqrt{\nu^\lambda}$  for large  $\lambda$ . Characterization of RMI behavior on the subdiffusion scale provides additional insights into the operation of the algorithm; we discuss this further in the next subsection. Denote by  $\varphi(\cdot)$  and  $\Phi(\cdot)$  the standard normal density and distribution functions, respectively. Let  $\Rightarrow$  denote convergence in distribution.

**THEOREM 2.** Consider the inverted- $V$  model in steady state under the RMI routing algorithm in the QED regime (C1–C2). Then, as  $\lambda \rightarrow \infty$ ,

$$(\hat{I}^\lambda, (\hat{I}_1^\lambda, \dots, \hat{I}_K^\lambda)1_{\{\hat{i}^\lambda > 0\}}) \Rightarrow (\hat{I}, (\hat{I}_1, \dots, \hat{I}_K)1_{\{\hat{i} > 0\}}), \quad (4)$$

where  $\hat{I}$  and  $(\hat{I}_1, \dots, \hat{I}_K)$  are independent;

$$\mathbb{P}[\hat{I} \leq 0] = \left( 1 + \delta \frac{\Phi(\delta)}{\varphi(\delta)} \right)^{-1};$$

$\mathbb{P}[\hat{I} > x | \hat{I} > 0] = \Phi(\delta - x)/\Phi(\delta)$ ,  $x \geq 0$ ;  $\mathbb{P}[\hat{I} \leq x | \hat{I} \leq 0] = e^{\delta x}$ ,  $x \leq 0$ ; and  $(\hat{I}_1, \dots, \hat{I}_K)$  is zero-mean multivariate normal, with  $\mathbb{E}[\hat{I}_i \hat{I}_j] = a_i 1_{\{i=j\}} - a_i a_j$ .

**PROOF.** See Mandelbaum et al. (2011).  $\square$

The contribution of the theorem to understanding system behavior under RMI is twofold. First, RMI achieves the same server fairness as LISF and IR (with appropriate weights) in the sense that idle servers are distributed across the pools according to the pools relative capacities ( $a_i$ 's). This is of significance because RMI requires less information for its operation than LISF and, unlike IR, RMI does not utilize information on pool capacities. Second, the “quality” of allocation of idle servers across the pools under RMI is revealed. Specifically, the number of idle servers in a pool deviates from  $(c_i^\lambda/c^\lambda)I^\lambda$  (a number determined by pools relative capacities) by a random quantity (normally distributed) of the order  $\sqrt[4]{\nu^\lambda}$ . In view of the fact that the number of idle servers in a pool is proportional to  $\sqrt{\nu^\lambda}$ , fluctuations of order  $\sqrt[4]{\nu^\lambda}$  are negligible when  $\lambda$  is not small.

**REMARK 7 (EQUAL SERVICE RATES).** When the service rates are equal across the server pools, i.e.,  $\mu_1 = \mu_2 = \dots = \mu_K$ , then  $\mu = \hat{\mu} = \mu_1$ ,  $\delta = \beta$ , and  $N^\lambda/\nu^\lambda \rightarrow 1$ , as  $\lambda \rightarrow \infty$ , and one recovers the well-known Erlang-C QED approximation (Halfin and Whitt 1981).

**REMARK 8.** Note that  $\sum_{i=1}^K \hat{I}_i^\lambda 1_{\{\hat{i}^\lambda > 0\}} = 0$  by definition. The limit (4) is consistent with this condition:  $\mathbb{P}[\sum_{i=1}^K \hat{I}_i = 0] = 1$ , because  $\mathbb{E}(\sum_{i=1}^K \hat{I}_i)^2 = 0$ .

**REMARK 9.** The dimensionality reduction (on the diffusion scale) can be deduced from the hydrodynamical equations in Dai and Tezcan (2011). For example, consider the case when there are  $K = 2$  pools. If the fraction of idle servers that are in the first pool exceeds  $c_1^\lambda/c^\lambda$ , then a disproportionate number of customers will be routed to the first pool, resulting in a lower number of idle servers in the first pool. As similar situation occurs when the fraction of idle servers that are in the first pool drops below  $c_1^\lambda/c^\lambda$ . Hence, the ratio of the number of idle servers in different pools should be equal to the ratio of the pool capacities.

**EXAMPLE 1 (PROBABILITY OF DELAY).** Consider an inverted- $V$  system with the following parameters:  $K = 2$ ,  $q_2 = 2q_1 = 2/3$ , and  $\mu_1 = 2\mu_2 = 2$  (e.g., patients/week). The total number of servers (e.g., beds) is taken to be

$$N^\lambda = \left\lceil \frac{\lambda + 0.5\sqrt{\lambda}}{q_1\mu_1 + q_2\mu_2} \right\rceil,$$

whereas  $N_1^\lambda = \text{round}(q_1 N^\lambda)$  and  $N_2^\lambda = N^\lambda - N_1^\lambda$ . We vary the arrival rate  $\lambda$  from 10 to 500—this corresponds to varying  $N_1^\lambda$  and  $N_2^\lambda$  from 3 to 128 and 6 to 256, respectively. In Figure 2, we plot the exact probability of delay along its QED approximation. Note that because of the PASTA property, the probability of delay is equal to  $\mathbb{P}[I^\lambda \leq 0]$ , and, hence, Theorem 2 renders

$$\mathbb{P}[\text{delay}] \rightarrow \left( 1 + \delta \frac{\Phi(\delta)}{\varphi(\delta)} \right)^{-1}, \quad (5)$$

as  $\lambda \rightarrow \infty$ . The preceding limit serves as the basis for calculating QED approximation for the finite system. To this end, the required QED parameter  $\delta$  changes slightly with  $\lambda$  because the number of servers in each pool is a natural number. Thus, when evaluating the QED approximation, we compute  $\delta$  for each value of  $\lambda$ :

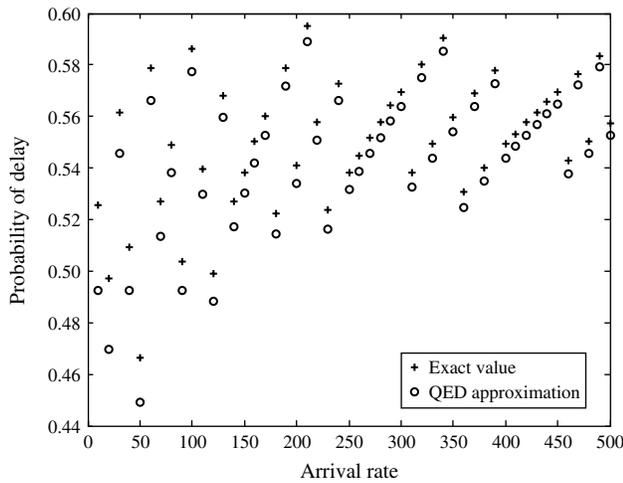
$$\delta^\lambda = \left( 1 - \frac{\lambda}{N_1^\lambda \mu_1 + N_2^\lambda \mu_2} \right) \sqrt{\lambda/\hat{\mu}^\lambda},$$

where

$$\hat{\mu}^\lambda = \frac{N_1^\lambda \mu_1^2 + N_2^\lambda \mu_2^2}{N_1^\lambda \mu_1 + N_2^\lambda \mu_2}.$$

A QED approximation for the probability of delay is obtained by evaluating the right-hand side of (5), with

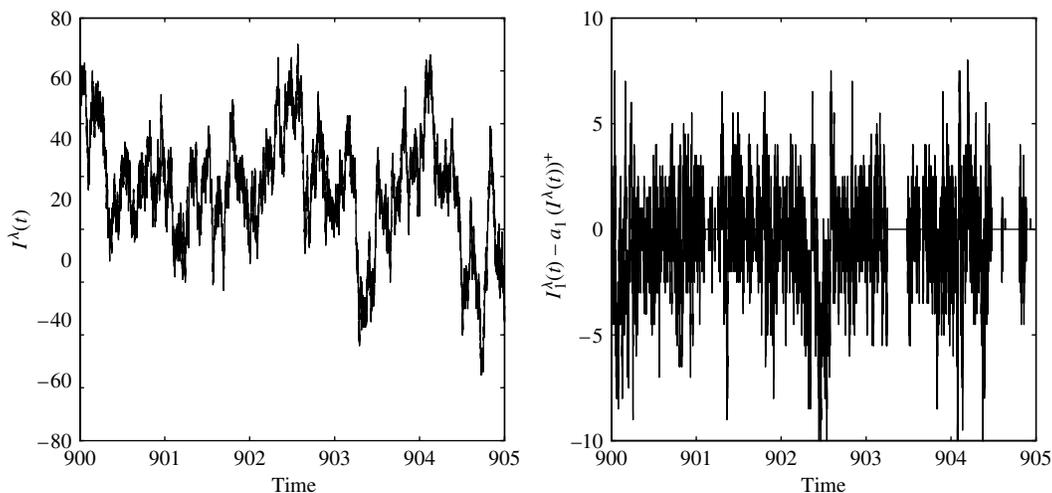
**Figure 2** Exact Values and QED Approximations of the Probability of Delay for the Sequence of Systems Described in Example 1



$\delta = \delta^\lambda$ . As can be seen in Figure 2, the QED approximation has a reasonable (useful) accuracy when system sizes and service rates are similar to those in Anonymous Hospital (see §2).

**EXAMPLE 2 (DIMENSIONALITY REDUCTION).** To illustrate the dimensionality reduction that arises in Theorem 2, we simulated an inverted- $V$  system under the RMI policy, with parameters that adhere to the QED scaling (C1) and (C2):  $K = 2$ ,  $\lambda = 3,950$ ,  $\mu_1 = 15$ ,  $\mu_2 = 7.5$ ,  $N_1 = 138$ , and  $N_2 = 276$  ( $a_1 = a_2 = 1/2$ ,  $\hat{\mu} = 11.25$ ,  $\delta \approx 0.86$ ). In Figure 3, we plot typical realizations of the total number of idle servers,  $\{I^\lambda(t), t \geq 0\}$ , and the centered number of idle server in the first pool,  $\{I_1^\lambda(t) - a_1(I^\lambda(t))^+, t \geq 0\}$ . Initially, at time  $t = 0$ , there are no idle servers and no customers await service, i.e.,  $I^\lambda(0) = 0$ . By time  $t = 900$ , the system is close to its stationary regime (there

**Figure 3** Illustration for Example 2



are more than  $7 \cdot 10^6$  arrivals/departures in the time interval  $[0, 900]$ ). One can observe that the processes  $\{I^\lambda(t), t \geq 0\}$  and  $\{I_1^\lambda(t) - a_1(I^\lambda(t))^+, t \geq 0\}$  evolve on two different counting scales—the first one on the  $\sqrt{\nu^\lambda}$ -scale ( $\sqrt{\nu^\lambda} \approx 18.7$ ), and the second one on the  $\sqrt[4]{\nu^\lambda}$ -scale ( $\sqrt[4]{\nu^\lambda} \approx 4.3$ ).

From the following corollary, it is immediate that, under RMI, the idleness ratios satisfy  $(1 - \rho_i^\lambda)/(1 - \rho_j^\lambda) \rightarrow \mu_i/\mu_j$  as  $\lambda \rightarrow \infty$ ; that is, the RMI policy achieves the same idleness ratios as the LIFS policy.

**COROLLARY 1.** Consider the inverted- $V$  model in steady state under the RMI routing algorithm in the QED regime (C1–C2). Then, as  $\lambda \rightarrow \infty$ ,  $\mathbb{E}(\hat{I}^\lambda)^- \rightarrow \mathbb{E}(\hat{I})^-$ ,  $\mathbb{E}(\hat{I}^\lambda)^+ \rightarrow \mathbb{E}(\hat{I})^+$ , and  $\mathbb{E}I_i^\lambda/\sqrt{\nu^\lambda} \rightarrow a_i\mathbb{E}(\hat{I})^+$  for  $i = 1, \dots, K$ , where  $\hat{I}$  is as in Theorem 2.

**PROOF.** See Mandelbaum et al. (2011).  $\square$

**REMARK 10 (LOSS OF PERFORMANCE).** As stated in the introduction, FSF routing is asymptotically optimal (within the set of all nonpreemptive, nonanticipating FCFS policies) in the sense that it stochastically minimizes the stationary queue length and waiting time as the arrival rate and number of servers grow large in the considered inverted- $V$  model. When analyzing the trade-off between fairness and performance (i.e., RMI versus FSF), one must consider the following two aspects.

On one hand, within our mathematical model, fairness comes at the cost of a decrease in system performance (e.g., an increase of the probability of delay). Theorem 2 (and Corollary 1), together with results

on FSF routing in Armony (2005), provides means for quantifying the loss of performance under the RMI policy (relative to FSF). In particular, results in Armony (2005) indicate that, under the FSF policy,

$$\mathbb{P}[\hat{I} \leq 0] = \left(1 + \delta_* \frac{\Phi(\delta_*)}{\varphi(\delta_*)}\right)^{-1};$$

$\mathbb{P}[\hat{I} > x | \hat{I} > 0] = \Phi(\delta_* - x\sqrt{\mu_\wedge/\hat{\mu}})/\Phi(\delta_*)$ ,  $x \geq 0$ ;  $\mathbb{P}[\hat{I} \leq 0] = e^{\delta x}$ ,  $x \leq 0$ ;  $\mu_\wedge = \min \mu_i$ ; and  $\delta_* = \delta\sqrt{\hat{\mu}/\mu_\wedge}$ . Thus, the probability of delay under the RMI policy is higher (because  $\delta_* \geq \delta$ ) than the probability of delay under the FSF policy by a factor of

$$\frac{1 + \delta_* \Phi(\delta_*)/\varphi(\delta_*)}{1 + \delta \Phi(\delta)/\varphi(\delta)}, \quad (6)$$

which depends on the spare capacity parameter  $\delta$  and the ratio of the arithmetic-mean service rate  $\hat{\mu}$  to the minimum service rate  $\mu_\wedge$ ; recall that in heavy traffic, only servers with the smallest service rate are idled under the FSF policy. The increase in the expected waiting time is given by the same ratio. This follows from the fact that the conditional expected wait, given that it is positive, is the same under the two routing policies. The decrease in performance should be carefully reviewed to determine if increased delays are clinically acceptable. For example, in Anonymous Hospital (see Table 1),  $\hat{\mu} \approx 1.18$  patients/week,  $\mu \approx 1.08$  patients/week, and  $\delta \approx 0.87$ , leading to the ratio in (6) being equal to 1.07, i.e., the probability of delay increases by 7% when employing the RMI instead of the FSF policy.

On the other hand, as stated in §1, the issue of fairness needs to be examined beyond the mathematical model with fixed service rates. In particular, ensuring fairness toward medical staff should be viewed within the context of providing a right set of incentives for staff. Although the waiting time of customers (patients) is stochastically minimized under the FSF policy, medical staff working in wards with nonminimal service rates have no incentive to further reduce LOS—any improvement would lead to a higher load. Thus, although the FSF policy nominally results in shorter waiting times, the RMI policy can be more beneficial for customers in the long run because the implementation of the RMI policy can eventually lead to shorter LOSs and, as a result, shorter waiting times.

*Weighted RMI.* Finally, we note that the RMI algorithm can be generalized to a weighted RMI (WRMI) algorithm as follows. Given a set of weights  $w_i > 0$  ( $i = 1, \dots, K$ ) such that  $\sum_{i=1}^K w_i = 1$ , the WRMI algorithm routes a customer to pool  $i$ , at time  $t$ , with probability  $I_i^{w,\lambda}(t-)/\sum_{j=1}^K I_j^{w,\lambda}(t-)$ , where  $I_i^{w,\lambda}(t) = w_i I_i^\lambda(t)$ . The weights can be used to adjust idleness ratios to a desired target (relative to the

ratio of service rates) in the QED regime. Unless all weights are equal, the resulting system is not reversible, and thus our analysis of RMI can not be extended to WRMI. However, insights gained from RMI can be used to heuristically analyze WRMI as follows. Consider a time instance  $t$  such that  $I^\lambda(t)/\sqrt{\nu^\lambda} > 0$ . Then, the departure rate of customers from pool  $i$  is  $c_i^\lambda - \mu_i I_i^\lambda(t)$  (where  $c_i^\lambda \gg \mu_i I_i^\lambda(t)$ ); on the other hand, customers enter service at pool  $i$  with rate  $\lambda I_i^{w,\lambda}(t)/\sum_{j=1}^K I_j^{w,\lambda}(t)$ . Therefore, for large  $\lambda$ ,

$$\frac{a_i}{a_j} \approx \frac{c_i^\lambda - \mu_i I_i^\lambda(t)}{c_j^\lambda - \mu_j I_j^\lambda(t)} \approx \frac{I_i^{w,\lambda}(t)}{I_j^{w,\lambda}(t)} = \frac{w_i I_i^\lambda(t)}{w_j I_j^\lambda(t)}.$$

In view of the preceding, one expects that the idleness ratios satisfy, as  $\lambda \rightarrow \infty$ ,

$$\frac{1 - \rho_i^\lambda}{1 - \rho_j^\lambda} \rightarrow \frac{w_j \mu_j}{w_i \mu_i}.$$

## 5.1. Comments on the Subdiffusion Scale

**5.1.1. Relevance and Implications.** Based on empirical data from our Anonymous Hospital, we estimated that  $\lambda \approx 189.7$  patients/week, and  $\hat{\mu} \approx 1.18$  patients/week (see Table 1); thus,  $\nu^\lambda = \lambda/\hat{\mu} \approx 160.8$ . These numbers and our theoretical analysis reveal that there exist three relevant counting scales (bed, room, subward) and three corresponding time scales (hour, day, week). To describe a stochastic process of interest (e.g., the number of available beds), one needs to define not only an appropriate counting scale, but also appropriate time intervals (scale) over which relevant changes in the process occur. On time intervals that are too short, no changes in the value of the process can be observed on the counting scale; on the other hand, on time intervals that are too long, the process covers the whole counting scale, and time variability can not be studied.

The finest counting scale is at the level of an individual bed/patient (order-1 scale), whereas the corresponding time scale is at the level of an hour (order-1/ $\lambda$ ). Indeed, when considering individual patient arrivals, the relevant time scale is based on hours because  $1/\lambda \approx 0.86$  hours. The coarsest counting scale (order- $\sqrt{\nu^\lambda}$ ; diffusion scale) is used to describe the total number of available beds and patients awaiting hospitalization. For a large hospital, this scale is defined by a subward because  $\sqrt{\nu^\lambda} \approx 12.7$  (beds or patients) is approximately the size of one-quarter to one-third of a ward, and the number of available beds/patients waiting is proportional to  $\sqrt{\nu^\lambda}$  (see Theorem 2). The corresponding time scale is defined by a “typical” LOS—a week in our case, because patients stay in a ward a bit less than a week on average (or  $1/\hat{\mu} \approx 0.85$  weeks; see Table 1).

The preceding two pairs of scales are standard for systems operating in the QED regime. The third pair of scales is intermediate and due to RMI. The counting scale (subdiffusion scale) is defined by a room (rooms in Anonymous Hospital have four beds) and is an order- $\sqrt[4]{\nu^\lambda}$  scale ( $\sqrt[4]{\nu^\lambda} \approx 3.6$  (beds or patients)). This scale is relevant in describing the number of patients that need to be moved between wards to make the *instantaneous* (at a specific moment of time) idleness ratios equal to the long-run (average) idleness ratios introduced earlier. The latter ratios do not provide information on the system behavior at a specific moment of time.

For example, consider the following two scenarios, assuming two symmetric wards: (i) the number of available beds in the two wards is the same at all times, and (ii) the number of available beds in the first ward is higher than in the second one for a long period of time, and then the situation is reversed for the same amount of time. Under both scenarios, the idleness ratio is unity. The room-level ( $\sqrt[4]{\nu^\lambda}$ ) counting scale provides information on deviations of the numbers of available beds (in steady state) from the numbers that ensure that the actual ratios of idle servers at a given moment of time are equal to the idleness ratios (average quantities). In particular, our results imply that fluctuations of instantaneous idleness ratios around their long-run averages are of the order  $1/\sqrt[4]{\nu^\lambda}$ . The associated time scale is based on a day (order- $1/\sqrt{\lambda\hat{\mu}}$ ;  $1/\sqrt{\lambda\hat{\mu}} \approx 0.87$  days) and describes relevant time intervals over which these relatively minor fluctuations of idleness ratios average out. Recall the two scenarios mentioned above. In the second one, several cycles are needed for idleness ratios to converge to unity. Indeed, if one observes the system for a short period of time only, then the system appears unfair. Yet, over long periods of time the system is fair. Finally, we note that the intermediate counting and time scales are related, in the sense that the larger the fluctuations of the number of available beds, the longer time intervals one requires for convergence of the idleness ratios.

Informally, our results indicate that, under the RMI policy, fluctuations of the numbers of idle beds in different wards, around values that ensure  $\mu_i/\mu_j$  instantaneous idleness ratios, are of the order  $\sqrt[4]{\nu^\lambda} \approx 3.6$  (beds or patients); i.e., the room-based counting scale is relevant. Indeed, Theorem 2 implies

$$I_i^\lambda \approx \frac{c_i^\lambda}{c^\lambda} I^\lambda + \hat{I}_i^\lambda \sqrt{I^\lambda},$$

and hence the idleness ratios obey

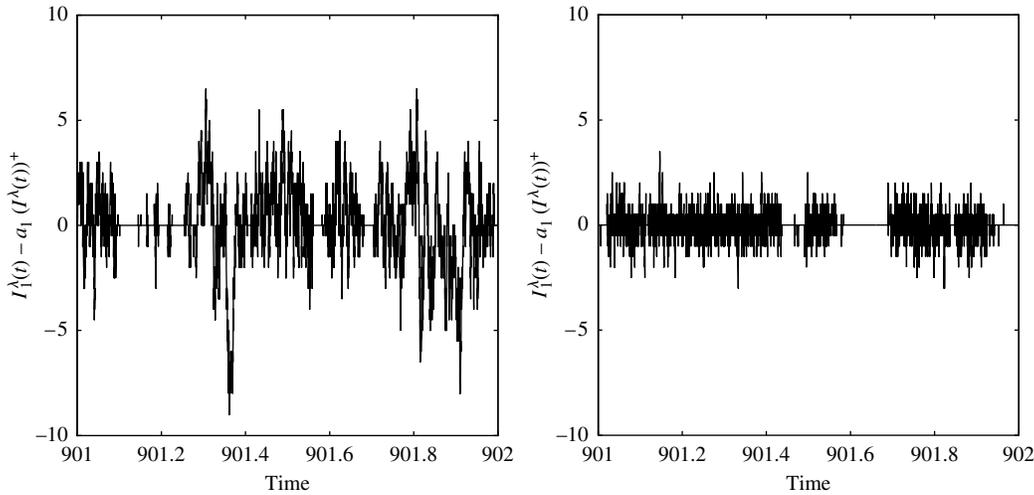
$$\frac{I_i^\lambda/N_i^\lambda}{I_j^\lambda/N_j^\lambda} \approx \frac{\mu_i}{\mu_j} \left( 1 + \frac{1}{\sqrt{I^\lambda}} \left( \hat{I}_i^\lambda/a_i - I_j^\lambda/a_j \right) \right).$$

Moreover, when available beds exist, even the minor differences (order- $1/\sqrt{I^\lambda}$  or, equivalently, order- $1/\sqrt[4]{\nu^\lambda}$ ) between instantaneous idleness ratios and the corresponding long-run idleness ratios ( $\mu_i/\mu_j$ ) are averaged out on time intervals of order  $1/\sqrt{\lambda\hat{\mu}}$  (time intervals that contain  $\sqrt{\nu^\lambda}$ -order arrivals/departures are of interest here, because central limit theorem deviations are of the order  $\sqrt[4]{\nu^\lambda}$  in that case; consequently, this corresponds to time intervals of length  $\sqrt{\nu^\lambda}/\lambda = 1/\sqrt{\lambda\hat{\mu}}$ ). In Anonymous Hospital, this corresponds roughly to days ( $1/\sqrt{\lambda\hat{\mu}} \approx 0.87$  days), and consequently one expects desired idleness ratios (with very high accuracy) on a weekly basis. Interestingly, the same behavior of intermediate scales occurs when the system operates under the LISF rule (see below).

**5.1.2. Technical Discussion.** Even though the three considered algorithms operate in different ways, they result in the same behavior on the diffusion scale. However, differences arise on the subdiffusion scale. Informally, Theorem 2 states that, for large  $\lambda$ , RMI deviations of  $I_i^\lambda$  around  $c_i^\lambda/c^\lambda (I^\lambda)^+$  are on the order of  $\sqrt[4]{\nu^\lambda}$ ; note that both  $I_i^\lambda$  and  $(I^\lambda)^+$  are  $\sqrt{\nu^\lambda}$ -order random variables. On the other hand, under IR policy,  $I_i^\lambda(t) - a_i (I^\lambda)^+$  is an order-1 random variable, as  $\lambda \rightarrow \infty$ . Consequently, although implementing RMI in a hospital setting will not ensure that an instantaneous idleness ratio is equal exactly to the desired value  $\mu_i/\mu_j$  at all times, the number of available beds in a ward will only differ from the desired one by a  $\sqrt[4]{\nu^\lambda}$ -quantity, which is negligible in comparison with the number of available beds in a ward.

Furthermore, it should be noted that, under both the IR and RMI policies, there exists a separation of time scales; namely, the processes  $\{I^\lambda(t)/\sqrt{\lambda}, t \geq 0\}$  and  $\{I_i^\lambda(t)/\sqrt{\lambda}, t \geq 0\}$  evolve on the order-1 time scale under both policies. This is typical in the QED regime—no time speedup is needed, in contrast to the case of conventional heavy traffic. However, the processes  $\{I_i^\lambda(t) - a_i (I^\lambda(t))^+, t \geq 0\}$  and  $\{(I_i^\lambda(t) - c_i^\lambda/c^\lambda (I^\lambda(t))^+)/\sqrt[4]{\nu^\lambda}, t \geq 0\}$  evolve on the  $\lambda^{-1}$ - and  $(\lambda\hat{\mu})^{-1/2}$ -time scales, as  $\lambda \rightarrow \infty$ , under IR and RMI routing, respectively. Hence, there is a separation of time scales, because the latter processes evolve on much faster time scales (order- $\lambda^{-1}$  and order- $(\lambda\hat{\mu})^{-1/2}$ ) than the former processes (order-1 time scale). In Figure 4, we plot typical sample paths of  $\{I_1^\lambda - a_1 (I^\lambda(t))^+, t \geq 0\}$  under RMI and IR routing, for the system described in Example 2—the difference in counting and time scales is evident. Therefore, in the context of Anonymous Hospital, the desired idleness ratios are maintained not only in the long run, but also on shorter time intervals. In particular, provided that idle servers exist ( $I^\lambda(0)\sqrt{\nu^\lambda} > \epsilon$  for some  $\epsilon > 0$ ), intervals of order  $\lambda^{-1}$  and  $(\lambda\hat{\mu})^{-1/2}$  are required under IR and RMI, respectively, for convergence of

Figure 4 Typical Sample Paths of  $\{I_1^\lambda(t) - a_1(I^\lambda(t))^+, t \geq 0\}$  in an Inverted- $V$  System Under the RMI (left) and IR (right) Policies



Note. Parameters of the system are the same as in Example 2.

the empirical idleness ratios to the long-run averages  $(1 - \rho_i^\lambda)/(1 - \rho_j^\lambda)$ , in the sense that (for large  $\lambda$ )

$$\sqrt{\nu^\lambda} \left( \frac{(1/N_i^\lambda) \int_0^{t/\lambda} I_i^\lambda(u) du}{(1/N_j^\lambda) \int_0^{t/\lambda} I_j^\lambda(u) du} - \frac{1 - \rho_i^\lambda}{1 - \rho_j^\lambda} \right) \text{ and}$$

$$\sqrt[4]{\nu^\lambda} \left( \frac{(1/N_i^\lambda) \int_0^{t/\sqrt{\lambda\hat{\mu}}} I_i^\lambda(u) du}{(1/N_j^\lambda) \int_0^{t/\sqrt{\lambda\hat{\mu}}} I_j^\lambda(u) du} - \frac{1 - \rho_i^\lambda}{1 - \rho_j^\lambda} \right)$$

converge to 0 as  $t$  increases ( $t$  is an order-1 quantity here) for IR and RMI, respectively.

In certain cases, the time scale separation can be used to explicitly evaluate the subdiffusion behavior of the system under IR routing in the QED regime. As seen in the following example, the idea is to exploit the fact that the subdiffusion process evolves on a faster time scale than diffusion processes. Recall that the subdiffusion behavior under RMI is characterized in Theorem 2.

EXAMPLE 3 (SUBDIFFUSION SCALE UNDER IR). Consider the inverted- $V$  model in the QED regime, under IR routing, with  $K = 2$ ,  $w_1 = w_2 = a_1 = a_2 = 1/2$ , and note that  $(I_1^\lambda(t) - I^\lambda(t)/2, I_2^\lambda(t) - I^\lambda(t)/2) = ((I_1^\lambda(t) - I_2^\lambda(t))/2, (I_2^\lambda(t) - I_1^\lambda(t))/2)$ . The heavy-traffic averaging principle (e.g., see Coffman et al. 1995) states that, when considering the distribution of  $(I_1^\lambda(t) - I_2^\lambda(t))$ , as  $\lambda \rightarrow \infty$ , one can act as if the total number of (scaled) idle servers is fixed (Whitt 2002, p. 70). In particular, on the event  $\{I^\lambda(t)/\sqrt{\nu^\lambda} > \epsilon\}$ ,  $\epsilon > 0$ , the distribution of  $(I_1^\lambda(t) - I_2^\lambda(t))$  converges, as  $\lambda \rightarrow \infty$ , to the stationary distribution of the birth–death continuous-time Markov chain, with transition rates  $r_{i,i+1} = 1/2 + 1_{\{i < 0\}} + (1 - \chi)1_{\{i=0\}}$  and  $r_{i,i-1} = 1/2 + 1_{\{i > 0\}} + \chi 1_{\{i=0\}}$  (here we assume that, if the two pools have the same number of idle servers, then a customer is routed to

the first pool with probability  $\chi \in [0, 1]$ ). Indeed, if  $(I_1^\lambda(t) - I_2^\lambda(t))$  is positive, then departures from pool 2 and new arrivals contribute to a decrease of this quantity; on the other hand, departures from pool 1 increase  $(I_1^\lambda(t) - I_2^\lambda(t))$ . This leads, for any  $\epsilon > 0$ , to

$$\mathbb{P}[I_1^\lambda(t) - I_2^\lambda(t) = i \mid I^\lambda(t)/\sqrt{\nu^\lambda} > \epsilon] \rightarrow (1 + 2(1 - \chi)1_{\{i > 0\}} + 2\chi 1_{\{i < 0\}})3^{-|i|-1},$$

as  $\lambda \rightarrow \infty$ .

We conjecture that the subdiffusion behavior of the system under the LISF algorithm is the same as the one under RMI. The conjecture is based on the following heuristic reasoning. A way to implement the LISF policy is to have servers completing service join a queue of idle servers. This queue operates in an FCFS fashion. Whenever a customer needs to be assigned to a server, it is routed to the server at the head of the queue. Observe that, under the described scheme, the server that has been idle for the longest time is assigned a customer before any other server. The state of the queue of idle servers at time  $t$  is an ordered list that consists of pool labels  $(1, 2, \dots, K)$ . Now, consider an inverted- $V$  model in the QED regime ( $\lambda \rightarrow \infty$ ) at a time instance  $t$  such that  $I^\lambda(t)/\sqrt{\nu^\lambda} > 0$ . Then, all the servers in the idle queue joined the queue within a time interval that is on the order of  $1/\sqrt{\lambda\hat{\mu}}$ ; i.e.,  $I^\lambda(t)/\sqrt{\nu^\lambda}$  remains approximately constant during this interval of time. The server pool labels in the idle queue are approximately independent, with a label being equal to  $i$  with probability

$$\frac{c_i^\lambda - \mu_i I_i^\lambda(t)}{c^\lambda - \sum_{i=1}^K \mu_i I_i^\lambda(t)} \approx \frac{c_i^\lambda - a_i \mu_i I^\lambda(t)}{c^\lambda - \hat{\mu} I^\lambda(t)} = a_i + \Theta(1/\sqrt{\nu^\lambda})$$

for large  $\lambda$ ; the standard asymptotic notation  $\Theta(1/\sqrt{\nu^\lambda})$  indicates that the second term is of the

order  $1/\sqrt{\nu^\lambda}$ . The preceding equation is due to the fact that a given label is of type  $i$  if a server in pool  $i$  completes service before any other server in the other pools. As a consequence, the random variable  $(I_1^\lambda(t) - a_1 I^\lambda(t))$  is of the order  $\sqrt[4]{\nu^\lambda}$ , because of the central limit theorem and the fact that the total number of labels in the queue is  $I^\lambda(t)$ , a quantity proportional to  $\sqrt{\nu^\lambda}$ .

## 6. Concluding Remarks

We considered routing algorithms that are applicable to routing hospital patients from the emergency department to internal wards. Given the heterogeneity of the wards, the objective is to achieve fairness from the point of view of hospital staff, while not hurting efficiency too severely. Wards are modeled as server pools, and two types of quantities are used to quantify fairness: *idleness* ratios and *flux* ratios. Under the LISF policy, which is considered to be “fair” and is commonly used in call centers, both ratios tend to the *ratios of service rates* in the respective server pools, when the system is in the QED regime; the appropriateness of the QED regime in modeling the ED-to-IW process is supported by empirical data, collected in Anonymous Hospital. In other words, LISF routing leads to a desirable outcome: faster servers work less, yet they produce more (serve more customers). However, the applicability of LISF routing in hospitals is limited because the algorithm requires information unavailable in hospitals on a real-time basis. The same idleness and flux ratios can be achieved by IR routing; yet, in that case, one must estimate (time-varying) ward (server pool) service capacities. We thus propose a randomized routing policy, RMI, that attains the same desired performance ratios as LISF, but requires only the number of idle servers in server pools for making decisions. The policy can be implemented in a hospital by using, for example, patient ID numbers as sources of randomness. A generalized version of RMI, WRMI, can be used to fine-tune the desired outcome.

The three algorithms (LISF, IR, and RMI) share one feature in common; namely, these algorithms achieve the desired idleness ratios (equal to service-rate ratios) by attempting to maintain the ratios of the numbers of idle servers in different pools equal to these idleness ratios at all times. However, because idleness ratios represent an average performance measure, one can vary the instantaneous ratios of idle servers depending on the total number of idle servers and still achieve the target (average) idleness ratios. This approach was proposed by Armony and Ward (2010). The advantage of such an algorithm is that it delivers a lower average waiting time compared to LISF, while maintaining the same long-run idleness ratios as LISF. However, the gain in performance does

not come for free. In particular, one must determine the optimal number of idle servers ratios as a function of the total number of servers and, therefore, the parameters of the system must be known (arrival rate, pool service capacities), i.e., the policy is effectively not blind.

### 6.1. Partial Information Routing—Simulation Analysis

We mentioned above that availability of information is critical for determining an appropriate routing policy in hospitals. Our proposed routing, RMI, requires the information on the number of available beds at each ward at the moment of routing (information that is quite minimal compared to the other routing policies considered in this paper). However, the occupancy status in the IWs is not available on a real-time basis in Anonymous Hospital; instead, the ED relies on one bed census update per day (in the morning). Thus, to implement RMI routing, it is necessary to estimate the system state at decision times, based on the system state at the last update time point.

An example of such partial-information routing can be found in Tseytlin and Zviran (2009), where the authors created a computer simulation model of the ED-to-IW process in Anonymous Hospital. They used it to examine various routing policies, according to some fairness and performance criteria, while accounting for the scarcity of information in the system. Simulating the ED-to-IW process helped achieve additional practical insights, by accommodating some analytically intractable features (such as time-inhomogeneous Poisson arrivals), and allowed analysis of more complex routing algorithms. The best-performing algorithm (in terms of both staff fairness and operational performance) proposed by Tseytlin and Zviran (2009) was an algorithm that minimized at each decision point a convex combination of the two conflicting demands: balanced occupancy rates and balanced flux. The implementation under partial information resulted in almost no deterioration of performance.

To conclude, we now explain briefly the way that the lacking information is predicted. Denote by  $M_j$  the number of occupied beds (busy servers) in ward (pool)  $j$ . This counter is updated at some time point  $T$  and is estimated at other decision time points according to the patient routing and the ward service rate; namely, we estimate the number of occupied beds in ward  $j$  at time  $t > T$  to be equal to  $(M_j - M_j \mu_j (t - T))^+$ ,  $j = 1, 2, 3, 4$ . After a routing decision is made, we update  $M_k = M_k + 1$  ( $k$  denotes the ward chosen to admit the next patient).

### 6.2. Routing at the Level of Individual Providers

Instead of examining fairness via our bed-based model, one can study fairness by means of a more

complex staff-based model. Such a more detailed model could potentially be used to study fairness at the level of individual care provides rather than wards. (It would thus be more relevant to the U.S. healthcare system, where typically nurses do not have fixed ward assignments and are paid on an hourly basis.) In that case, one still needs to keep track of the number of patients in wards. However, two additional aspects must be modeled explicitly: (i) patient “service” requests, i.e., when certain tasks need to be performed by nurses/doctors (for example, in Belgium, nurses work is broken down into 23 representative tasks for planning purposes; Williams 2000); and (ii) the process of assigning of these tasks to individual staff members. For example, a model can be constructed from the following components:

- Parallel-server systems with heterogeneous servers (Atar 2008): A system represents a ward, and servers represent medical staff; tasks in a ward are assigned to individual doctors/nurses according to a (ward) scheduling policy.

- Erlang-R (“R” for “reentrant”) model with a bounded number of customers (Yom-Tov and Mandelbaum 2011): The reentrant aspect of the model captures the fact that customers (patients) require service (tasks) multiple times during their sojourn times (LOS) in the system; the bounded number of customers is due to the finite number of beds in a ward.

- A (hospital) routing policy that assigns customers (patients) to one of parallel Erlang-R systems (wards) with heterogeneous servers.

Overall fairness can be achieved by enforcing fairness both at the inter- and intraward levels. There exist multiple time scales in this model: a task scale (minutes/tens of minutes), a “content” scale (tens of minutes/hours; time between tasks for a single patient), a shift scale (hours), and finally a length-of-stay (a week) scale. These scales can potentially be used to simplify the analysis of the system. Indeed, during a patient’s stay in the hospital, not only many tasks are performed, but also many shifts are rotated. Hence, one expects that a patient experiences an “equivalent” service (task) rate on the LOS scale. This equivalent rate is a function of staff in the ward as well as the task assignment policy. As a consequence, when considering routing at the hospital level, one can plausibly replace the heterogeneous Erlang-R model with a parallel-server system with an equivalent service rate, i.e., one obtains our inverted- $V$  model.

### 6.3. Future Research

We propose a number of research directions motivated by our work. First, robustness of our insights against distributional assumptions on service times remains to be investigated. Second, the inverted- $V$  model takes into account primarily the number of beds

in each ward, whereas hospital staff (nurses and doctors) affect the model indirectly through server service rates. The model can be improved by explicitly modeling staff as well. In that case, two-scale (doctors and beds) models arise. Third, patients to be hospitalized in the IWs are classified into several categories. When arriving to the ED, patients are classified as “walking” or “lying”; in addition, prior to running the Justice Table, they are classified as “regular,” “special care,” or “ventilated.” The load, imposed on the hospital by patients, varies significantly among different categories, in LOS, complexity of treatment, and waiting times. Thus, it is of prime interest to extend our model to accommodate multiple customer classes. Last, in the present work, service rates are taken to be exogenous quantities, i.e., there is no attempt to capture possible dependency between the routing algorithm and service rates of doctors and nurses. However, such dependency does exist because the hospital staff adapts to routing policies by increasing/decreasing their service rates and/or quality of care. Tools from game theory can be applicable in modeling such effects.

### Acknowledgments

Avishai Mandelbaum’s research was supported in part by the Binational Science Foundation [Grants 2005175 and 2008480]; the Israeli Science Foundation [Grant 1357/08]; and by the Technion funds for the promotion of research and sponsored research. Petar Momčilović’s research was supported in part by the National Science Foundation [Grant CNS-0643213].

### References

- Armony, M. 2005. Dynamic routing in large-scale service systems with heterogeneous servers. *Queueing Syst. Theory Appl.* 51(3–4) 287–329.
- Armony, M., A. Mandelbaum. 2012. Routing and staffing in large-scale service systems: The case of homogeneous impatient customers and heterogeneous servers. *Oper. Res.* 59(1) 50–65.
- Armony, M., A. Ward. 2010. Fair dynamic routing policies in large-scale systems with heterogeneous servers. *Oper. Res.* 58(3) 624–637.
- Armony, M., S. Israelit, A. Mandelbaum, Y. Marmor, Y. Tseytlin, G. Yom-Tov. 2011. Patient flow in hospitals: A data-based queueing-science perspective. Working paper, New York University, New York.
- Atar, R. 2008. Central limit theorem for a many-server queue with random service rates. *Ann. Appl. Probab.* 18(4) 1548–1568.
- Atar, R., Y. Y. Shaki, A. Shwartz. 2011. A blind policy for equalizing cumulative idleness. *Queueing Systems Theory Appl.* 67(4) 275–293.
- Avi-Itzhak, B., H. Levy. 2004. On measuring fairness in queues. *Adv. Appl. Probab.* 36(3) 919–936.
- Bekker, R., A. M. de Bruin. 2010. Time-dependent analysis for refused admissions in clinical wards. *Ann. Oper. Res.* 178(1) 45–65.
- Ben-Zrihen, M., J. Borsher, A. Reiss, Y. Tseytlin. 2007. Behavioral models in customer service centers. IE&M project, Industrial Engineering and Management, Technion, Haifa, Israel.
- Borst, S., A. Mandelbaum, M. Reiman. 2004. Dimensioning of large call centers. *Oper. Res.* 52(1) 17–34.

- Cabral, F. B. 2005. The slow server problem for uninformed customers. *Queueing Syst. Theory Appl.* **50**(4) 353–370.
- Cabral, F. B. 2007. Queues with heterogeneous servers and uninformed customers: Who works the most? Working paper, Universidade Federal de Pelotas, UNIPAMPA, Brazil.
- Coffman, E. G., Jr., A. A. Puhalskii, M. I. Reiman. 1995. Polling systems with zero switchover times: A heavy-traffic averaging principle. *Ann. Appl. Probab.* **5**(3) 681–719.
- Colquitt, J., D. Conlon, M. Wesson, C. Porter, K. Y. Ng. 2001. Justice at the millennium: A meta-analytic review of 25 years of organizational justice research. *J. Appl. Psych.* **86**(3) 425–445.
- Dai, J., T. Tezcan. 2011. State space collapse in many-server diffusion limits of parallel server systems. *Math. Oper. Res.* **36**(2) 271–320.
- de Bruin, A. M., R. Bekker, L. van Zanten, G. M. Koole. 2010. Dimensioning hospital wards using the Erlang loss model. *Ann. Oper. Res.* (1) 23–43.
- de Véricourt, F., O. B. Jennings. 2011. Nurse staffing in medical units: A queueing perspective. *Oper. Res.* **59**(6) 1320–1331.
- Elkin, K., N. Rozenberg. 2007. Patients' flow from the emergency department to the internal wards. IE&M project, Industrial Engineering and Management, Technion, Haifa, Israel.
- Erlang, A. K. 1948. On the rational determination of the number of circuits. E. Brockmeyer, H. L. Halstrom, A. Jensen, eds. *The life and works of A.K. Erlang*. Copenhagen Telephone Company, Copenhagen, 216–221.
- Gans, N., G. Koole, A. Mandelbaum. 2003. Telephone call centers: Tutorial, review, and research prospects. *Manufacturing Service Oper. Management* **5**(2) 79–141.
- González, P., C. Herrero. 2004. Optimal sharing of surgical costs in the presence of queues. *Math. Methods Oper. Res.* **59**(3) 435–446.
- Green, L. V. 2004. Capacity planning and management in hospitals. M. L. Brandeau, F. Sainfort, W. P. Pierskalla, eds. *Operations Research and Health Care: A Handbook of Methods and Applications*. International Series in Operations Research & Management Science. Kluwer, Norwell, MA, 15–42.
- Green, L.V. 2008. Using operations research to reduce delays for healthcare. Z.-L. Chen, S. Raghavan, eds. *Tutorials in Operations Research*. INFORMS, Hanover, MD, 1–16.
- Gurvich, I., W. Whitt. 2009. Queue-and-idleness-ratio controls in many-server service systems. *Math. Oper. Res.* **34**(2) 363–396.
- Halfin, S., W. Whitt. 1981. Heavy-traffic limits for queues with many exponential servers. *Oper. Res.* **29**(3) 567–588.
- Huseman, R. C., J. D. Hatfield, E. W. Miles. 1987. A new perspective on equity theory: The equity sensitivity construct. *Acad. Management Rev.* **12**(2) 222–234.
- Joint Commission on the Accreditation of Healthcare Organizations (JCAHO). 2004. JCAHO requirement: New leadership standard on managing patient flow for hospitals. *Joint Commission Perspectives* **24**(2) 13–14.
- Kc, D. S., C. Terwiesch. 2009. Impact of workload on service time and patient safety: An econometric analysis of hospital operations. *Management Sci.* **55**(9) 1486–1498.
- Kelly, F. 1979. *Reversibility and Stochastic Networks*. Wiley, New York.
- Larsen, R. L., A. K. Agrawala. 1983. Control of a heterogeneous two-server exponential queueing system. *IEEE Trans. Software Engrg.* **9**(4) 522–526.
- Larson, R. C. 1987. Perspectives on queues: Social justice and the psychology of queueing. *Oper. Res.* **35**(6) 895–905.
- Lin, W., P. R. Kumar. 1984. Optimal control of a queueing system with two heterogeneous servers. *IEEE Trans. Automat. Control* **29**(8) 696–703.
- Litvak, E., M. C. Long, A. B. Cooper, M. L. McManus. 2001. Emergency department diversion: Causes and solutions. *Acad. Emerg. Med.* **8**(11) 1108–1110.
- Mandelbaum, A., P. Momčilović, Y. Tseytlin. 2011. Online appendices to “On fair routing from emergency departments to hospital wards: QED queues with heterogeneous servers.” <http://iew3.technion.ac.il/serveng/References/>.
- McManus, M. L., M. C. Long, A. Cooper, E. Litvak. 2004. Queueing theory accurately models the need for critical care resources. *Anesthesiology* **100**(5) 1271–1276.
- New York Times*. 2002. 1 in 3 hospitals say they divert ambulances. (April 9), <http://www.nytimes.com/2002/04/09/us/1-in-3-hospitals-say-they-divert-ambulances.html>.
- Rafaeli, A., G. Barron, K. Haber. 2002. The effects of queue structure on attitudes. *J. Service Res.* **5**(2) 125–139.
- Ramakrishnan, M., D. Sier, P. G. Taylor. 2005. A two-time-scale model for hospital patient flow. *IMA J. Management Math.* **16**(3) 197–215.
- Rubinovitch, M. 1985a. The slow server problem. *J. Appl. Probab.* **22**(1) 205–213.
- Rubinovitch, M. 1985b. The slow server problem: A queue with stalling. *J. Appl. Probab.* **22**(4) 309–317.
- Stockbridge, R. H. 1991. A martingale approach to the slow server problem. *J. Appl. Probab.* **28**(2) 480–486.
- Tseytlin, Y. 2009. Queueing systems with heterogeneous servers: On fair routing of patients in emergency departments. M.Sc. thesis, Industrial Engineering and Management, Technion, Haifa, Israel.
- Tseytlin, Y., A. Zviran. 2009. Simulation of patients routing from an emergency department to internal wards in Anonymous Hospital. IE&M project, Industrial Engineering and Management, Technion, Haifa, Israel.
- Whitt, W. 2002. *Stochastic-Process Limits: An Introduction to Stochastic-Process Limits and Their Application to Queues*. Springer, New York.
- Williams, R. 2000. It all adds up. *Nurs. Stand.* **14**(31) 12–13.
- Wright, J., R. King. 2006. *We All Fall Down: Goldratt's Theory of Constraints for Healthcare Systems*. North River Press, Great Barrington, MA.
- Yom-Tov, G. 2007. Queues in hospitals: Semi-open queueing networks in the QED regime. Technical report, Industrial Engineering and Management, Technion, Haifa, Israel.
- Yom-Tov, G. 2010. Queues in hospitals: Queueing networks with ReEntering customers in the QED regime (QED = Quality- and Efficiency-Driven). Ph.D. thesis, Technion, Haifa, Israel.
- Yom-Tov, G., A. Mandelbaum. 2011. Erlang-R: A time-varying queue with ReEntrant customers, in support of healthcare staffing. Working paper, Industrial Engineering and Management, Technion, Haifa, Israel.
- Zhang, B., J. S. H. van Leeuwen, B. Zwart. 2012. Staffing call centers with impatient customers: Refinements to many-server asymptotics. *Oper. Res.* **60**(2) 461–474.