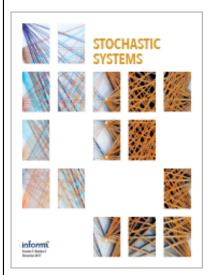
This article was downloaded by: [132.68.161.185] On: 12 November 2018, At: 02:49 Publisher: Institute for Operations Research and the Management Sciences (INFORMS) INFORMS is located in Maryland, USA



Stochastic Systems

Publication details, including instructions for authors and subscription information: http://pubsonline.informs.org

On Patient Flow in Hospitals: A Data-Based Queueing-Science Perspective

Mor Armony, Shlomo Israelit, Avishai Mandelbaum, Yariv N. Marmor, Yulia Tseytlin, Galit B. Yom-Toy

To cite this article:

Mor Armony, Shlomo Israelit, Avishai Mandelbaum, Yariv N. Marmor, Yulia Tseytlin, Galit B. Yom-Tov (2015) On Patient Flow in Hospitals: A Data-Based Queueing-Science Perspective. Stochastic Systems 5(1):146-194. https://doi.org/10.1287/14-SSY153

Full terms and conditions of use: http://pubsonline.informs.org/page/terms-and-conditions

This article may be used only for the purposes of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval, unless otherwise noted. For more information, contact permissions@informs.org.

The Publisher does not warrant or guarantee the article's accuracy, completeness, merchantability, fitness for a particular purpose, or non-infringement. Descriptions of, or references to, products or publications, or inclusion of an advertisement in this article, neither constitutes nor implies a guarantee, endorsement, or support of claims made of that product, publication, or service.

Copyright © 2015, The author(s)

Please scroll down for article—it is on subsequent pages

INFORMS is the largest professional society in the world for professionals in the fields of operations research, management science, and analytics.

For more information on INFORMS, its publications, membership, or meetings visit http://www.informs.org

ON PATIENT FLOW IN HOSPITALS: A DATA-BASED QUEUEING-SCIENCE PERSPECTIVE

By Mor Armony*, Shlomo Israelit[†], Avishai Mandelbaum[‡], Yariv N. Marmor[§], Yulia Tseytlin[¶], and Galit B. Yom-Tov^{||}

 NYU^* , Rambam Hospital[†], Technion[‡], ORT Braude College & Mayo Clinic[§], IBM Research[¶], Technion^{||}

Hospitals are complex systems with essential societal benefits and huge mounting costs. These costs are exacerbated by inefficiencies in hospital processes, which are often manifested by congestion and long delays in patient care. Thus, a queueing-network view of patient flow in hospitals is natural for studying and improving its performance. The goal of our research is to explore patient flow data through the lens of a queueing scientist. The means is exploratory data analysis (EDA) in a large Israeli hospital, which reveals important features that are not readily explainable by existing models.

Questions raised by our EDA include: Can a simple (parsimonious) queueing model usefully capture the complex operational reality of the Emergency Department (ED)? What time scales and operational regimes are relevant for modeling patient length of stay in the Internal Wards (IWs)? How do protocols of patient transfer between the ED and the IWs influence patient delay, workload division and fairness? EDA also underscores the importance of an integrative view of hospital units by, for example, relating ED bottlenecks to IW physician protocols. The significance of such questions and our related findings raise the need for novel queueing models and theory, which we present here as research opportunities.

Hospital data, and specifically patient flow data at the level of the individual patient, is increasingly collected but is typically confidential and/or proprietary. We have been fortunate to partner with a hospital that allowed us to open up its data for everyone to access. This enables reproducibility of our findings, through a user-friendly platform that is accessible via the Technion SEELab.

1. Introduction. Health care systems in general, and hospitals in particular, are major determinants of our quality of life. They also require a significant fraction of our resources and, at the same time, they suffer from (quoting a physician research partner) "a ridiculous number of inefficiencies;

Received June 2014.

Keywords and phrases: Queueing models, queueing networks, healthcare, patient flow, EDA, emergency departments, hospital wards, event logs.

thus everybody—patients, families, nurses, doctors and administrators are frustrated." In (too) many instances, this frustration is caused and exacerbated by delays—"waiting for something to happen"; in turn, these delays and the corresponding queues signal inefficiencies. Hospitals thus present a propitious ground for research in Queueing Theory and, more generally, Applied Probability (AP), Operations Research (OR) and Service Engineering (SE). Such research would ideally culminate in reduced congestion (crowding) and its accompanying important benefits: clinical, financial, psychological and societal. For such benefits to accrue, it is critical that the supporting research is data-based.

As it happens, however, operational hospital data is accessible to very few researchers, and patient-level data has in fact been publicly unavailable. The reasons span nonexistence or poor quality of data, concerns for patient confidentiality, and proprietorial constraints or lack of incentives for data owners. We attempt to address these issues as follows. First, we present an Exploratory Data Analysis (EDA) of a 1000-bed hospital, covering 3 years of patient-flow at the inter-departmental level of the individual patient. Through this EDA, we identify and propose research opportunities for AP, OR and SE. Then, we open up our operational database and make it universally accessible at the Technion IE&M Laboratory for Service Enterprise Engineering (SEELab): it can be either downloaded, or analyzed online with a user-friendly platform for EDA. Our goal is thus to provide an entry to and accelerate the learning of data-based OR of hospitals; researchers can use it (and some already have) to reproduce our EDA, which would serve as a trigger and a starting point for further data mining and novel research of their own.

1.1. Patient flow focus. Of particular interest to both researchers and practitioners is patient flow in hospitals: improving it can have a significant impact on quality of care as well as on patient satisfaction; restricting attention to it adds a necessary focus to our work. Indeed, the medical community has acknowledged the importance of patient flow management (e.g. Standard LD.3.10.10, which the Joint Commission on Accreditation of Hospital Organizations (JCAHO, 2004) set for patient flow leadership). This acknowledgment is natural, given that operational measures of patient flow are relatively easy to track, and that they inherently serve as proxies for other quality of care measures (see Section 6.1). In parallel, patient flow has caught the attention of researchers in OR in general, and Queueing Theory in particular. This is not surprising: hospital systems, being congestion-prone, naturally fit the framework of Queueing Theory, which captures the tradeoffs between (operational) service quality and resource efficiency.

Our starting point is that a queueing network encapsulates the operational dimensions of patient flow in a hospital, with the medical units being the nodes of the network; patients are the customers, while beds, medical staff and medical equipment are the servers. What are the special features of this queueing network? To address this question, we study an extensive data set of patient flow through the lens of a queueing scientist. Our study highlights interesting phenomena that arise in the data, which leads to a discussion of their implications on system operations and queueing modeling, and culminates in the proposal of related research opportunities.

However, patient flow is still too broad a subject for a single study. We thus focus on the inter-ward resolution, as presented in the flow chart (process map) of Figure 1; this is in contrast to intra-ward or out-of-hospital patient flow. We further narrow the scope to the relatively isolated ED+IW network, as depicted in Figure 2 and elaborated on in §1.2.1.

1.2. Rambam medical center. Our data originates from Rambam Medical Center, which is a large Israeli academic hospital. This hospital caters to a population of more than two million people, and it serves as a tertiary referral center for twelve district hospitals. The hospital consists of about 1000 beds and 45 medical units, with about 75,000 patients hospitalized annually. The data includes detailed information on patient flow throughout the hospital, over a period of several years (2004–2008), at the flow level of Figure 1, and the resolution level of individual patients. Thus, the data allows one to follow the paths of individual patients throughout their stay at the hospital, including admission, discharge, and transfers between hospital units.

1.2.1. The ED+IW network. Traditionally, hospital studies have focused on individual units, in isolation from the rest of the hospital; but this approach ignores interactions among units. On the flip side, looking at the hospital as a whole is complex and may lack necessary focus. Instead, and although our data encompasses the entire hospital, we focus on a sub-network that consists of the main Emergency Department (ED) (adult Internal, Orthopedics, Surgery, and Trauma) and five Internal Wards (IWs), denoted by A through E; see Figure 2. This sub-network, referred to as ED+IW, is more amenable to analysis than studying the entire hospital. At the same time, it is truly a system of networked units, which requires an integrative approach for its study. Moreover, the ED+IW network is also not too small: According to our data, approximately 47% of the patients entering the hospital remain within this sub-network, and 16% of those are hospitalized in the IWs. Finally, the network is fairly isolated in the sense that its interactions with the

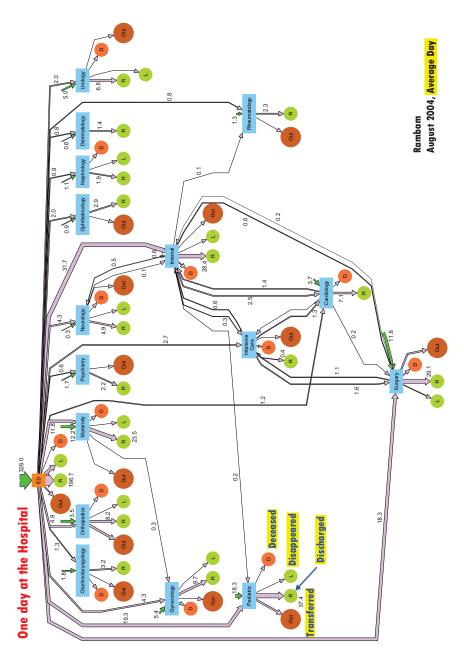
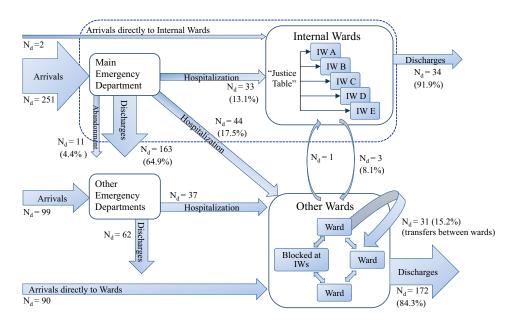


Fig 1. Patient Flow (Process Map) at inter-ward resolution. (Data animation is available at SEEnimations). For example, during the period over which the flow was calculated (August 2004), 326 patients arrived to the ED per day on average, and 18.3 transferred from the ED to Surgery. (To avoid clutter, arcs with monthly flow below 4 patients were filtered out; Created by SEEGraph, at the Technion SEELab.)



 N_d - daily average number of patients per weekdays (excluding holidays) total over 105 days, for period January 1, 2007 - May, 31, 2007

main ED - Internal, Surgery, Traumatology, and Orthopedic EDs

Fig 2. Patient flow in Rambam—zooming in on the ED+IW network.

rest of the hospital are minimal. To wit, virtually all arrivals into the ED are from outside the hospital, and 91.6% of the patient transfers into the IWs are either from outside the hospital or from within the ED+IW network.

1.2.2. Data description. Rambam's 2004–2008 patient-level flow data consists of 4 compatible "tables", that capture hospital operations as follows. The first table (Visits) contains records of ED patients, including their ID, arrival and departure times, arrival mode (e.g. independently or by ambulance), cause of arrival, and some demographic data. The second table (Justice Table) contains details of the patients that were transferred from the ED to the IWs. This includes information on the time of assignment from the ED to an IW, the identity of this IW, as well as assignment cancelations and reassignment times when relevant. The third table (Hospital Transfers) consists of patient-level records of arrivals to and departures from hospital wards. It also contains data on the ward responsible for each patient as, sometimes, due to lack of capacity, patients are not treated in the ward that is clinically most suitable for them; hence, there could be a distinction

between the physical location of a patient and the ward that is clinically in charge of that patient. The last table (Treatment) contains individual records of first treatment time in the IWs. Altogether, our data consists of over one million records.

1.3. "Apologies" to the statistician. Our approach of learning from data is in the spirit of Tukey's Exploratory Data Analysis (EDA) (Tukey, 1977), which gives rise to the following two "apologies". Firstly, the goals of the present study, as well as its target audience and space considerations, render secondary the role of "rigorous" statistical analysis (e.g. hypothesis testing, confidence intervals, model selection).

Secondly, our data originates from a single Israeli hospital, operating during 2004–2008. This casts doubts on the scope of the scientific and practical relevance of the present findings, and rightly so. Nevertheless, other studies of hospitals in Israel (Marmor (2003); Tseytlin (2009) and Section 5.6 of EV) and in Singapore (Shi et al., 2013), together with privately-communicated empirical research by colleagues, reveal phenomena that are common across hospitals worldwide (e.g. the LOS distributions in Figure 9). Moreover, our study can serve as a benchmark to compare against other hospitals. Finally, the present research has already provided the empirical foundation for several graduate theses, each culminating in one or several data-based theoretical papers (see §2.1).

- 1.4. Paper structure. The rest of the paper is organized as follows: We start with a short literature review in Section 2. We then proceed to discuss the gate to the hospital—the ED—in Section 3, followed by the IWs (§4), and the ED+IW network as a whole (§5). We start each section with background information. Next, we highlight relevant EDA, and lastly we propose corresponding research opportunities. In §6, we offer a final commentary, where we also provide a broader discussion of some common themes that arise throughout the paper. Finally, the Appendix covers data access instructions and documentation, as well as EDA logistics. We encourage interested readers to refer to EV: an online extended version of the present paper, which provides a more elaborate discussion of various issues raised here, and covers additional topics that we do not include due to focus and space considerations.
- 2. Some hints to the literature. Patient flow in hospitals has been studied extensively. Readers are referred to the papers in Hall (2013) and Denton (2013) which provide further leads to many other references. In the present subsection, we merely touch on published work, along the three

dimensions that are most relevant to our study: a network view, queueing models and data-based analysis. Many additional references to recent and ongoing research on particular issues that arise throughout the paper, will be cited as we go along. This subsection concludes with what can be viewed as a "proof of concept": a description of some existing research that the present work and our empirical foundation have already triggered and supported.

Most research on patient flow has concentrated on the ED and how to improve the internal ED flows. There are a few exceptions that offer a broader view. For example, Cooper et al. (2001) identifies a main source of ED congestion to be controlled variability, downstream from the ED (e.g. operatingroom schedules). In the same spirit, de Bruin et al. (2007) observes that "refused admissions at the First Cardiac Aid are primarily caused by unavailability of beds downstream the care chain." These blocked admissions can be controlled via proper bed allocation along the care chain of Cardiac inpatients; to support such allocations, a queueing network model was proposed, with parameters that were estimated from hospital data. Broadening the view further, Hall et al. (2006) develops data-based descriptions of hospital flows, starting at the highest unit-level (yearly view) down to specific sub-wards (e.g. imaging). The resulting flow charts are supplemented with descriptions of various factors that cause delays in hospitals, and then some means that hospitals employ to alleviate these delays. Finally, Shi et al. (2014) develops data-based models that lead to managerial insights on the ED-to-Ward transfer process.

There has been a growing body of research that treats operational problems in hospitals with Operations Research (OR) techniques. Brandeau, Sainfort and Pierskalla (2004) is a handbook of OR methods and applications in health care; the part that is most relevant to this paper is its chapter on Health Care Operations Management (OM). Next, Green (2008) surveys the potential of OR in helping reduce hospital delays, with an emphasis on queueing models. A recent handbook on System Scheduling is Hall (2012) which contains additional leads on OR/OM and queueing perspectives of patient flow. Of special interest is Chapter 8, where Hall describes the challenging reality of bed management in hospitals. Jennings and de Véricourt (2008, 2011) and Green and Yankovic (2011) apply queueing models to determine the number of nurses needed in a medical ward. Green (2004) and de Bruin et al. (2009) rely on queueing models such as Erlang-C and loss systems, to recommend bed allocation strategies for hospital wards. Lastly, Green, Kolesar and Whitt (2007) survey and develop (time-varying) queueing networks that help determine the number of physicians and nurses required in an ED.

There is also an increased awareness of the significant role that data can, and often must, play in patient flow research. For example, Kc and Terwiesch (2009) is an empirical work in the context of ICU patient flow; it has inspired the analytical model of Chan, Yom-Tov and Escobar (2014) (see also Chan, Farias and Escobar (2014) on the correlation between patient wait and ICU LOS). Another example is Baron et al. (2014) that does both modeling and data analysis for patient flow in outpatient test provision centers. More on patient flow in outpatient clinics and the need for relevant data is discussed in Froehle and Magazine (2013).

- 2.1. A proof of concept. The present research has provided the empirical foundation for several graduate theses and subsequent research papers: Marmor (2010) studied ED architectures and staffing (see Zeltyn et al. (2011) and Marmor et al. (2012)); Yom-Tov (2010) focused on time-varying models with reentrant customers in the ED (Yom-Tov and Mandelbaum, 2014) and the IWs; Tseytlin (2009) investigated the transfer process from the ED to the IWs (Mandelbaum, Momcilovic and Tseytlin, 2012); Maman (2009) explored over-dispersion characteristics of the arrival process into the ED (Maman, Zeltyn and Mandelbaum, 2011); and Huang (2013) developed scheduling controls that help ED physicians choose between newly-arriving vs. inprocess patients, while still adhering to triage constraints (Huang, Carmeli and Mandelbaum, 2015).
- 3. Emergency Department. The Emergency Department (ED) is the gate to the hospital, through which virtually all non-elective patients enter. Patient flow within the ED has been widely investigated, both academically (Hall et al., 2006; Saghafian, Austin and Traub, 2014; Zeltyn et al., 2011) and in practice (IHI, 2011; McHugh et al., 2011). Here we thus content ourselves with its empirical macro (black box) view. Specifically, we highlight interesting phenomena that relate to patient arrivals, departures, and occupancy counts. Our EDA underscores the importance of including time- and state-dependent effects in the ED—some of these are not readily explained by existing queueing models. Yet, our EDA also reveals that a simple stationary model may provide a good fit for patient-counts during periods when the ED is most congested. For limited purposes, therefore, our EDA supports the use of simple stationary models for the ED, which has been prevalent in the literature (e.g. de Bruin et al. (2009); Dong and Whitt (2014); Green et al. (2006)).
- 3.1. Basic facts. The main ED has 40 beds and it treats on average 251 patients daily: close to 60% are classified as Internal (general) patients

and the rest are Surgical, Orthopedic, or Multiple Trauma. While there are formally 40 beds in the ED, this bed capacity is highly flexible and can be doubled and more. Indeed, there is effectively no upper bound on how many patients can simultaneously reside within the ED—either in beds or stretchers, chairs, etc. The hospital has other EDs, physically detached from the main one discussed here—these are dedicated to other patient types such as Pediatrics or Ophthalmology. Throughout the rest of our paper, we focus on the main ED and simply refer to it as the ED. Furthermore, within the ED, we focus on Internal (general) patients, in beds or walking: they constitute the majority of ED patients and give rise to ample operational challenges.

During weekdays, the average length of stay (ALOS) of patients in the ED is 4.25 hours: this covers the duration from entry until the decision to discharge or hospitalize; it does not include boarding time, which is the duration between hospitalization decision to actual transfer. We estimate boarding time to be 3.2 hours on average (See Section 5.2). In addition, 10% (5%) of weekday patients experience LOS that is over 8 (11) hours, and about 3–5% leave on their own (LWBS = left without being seen by a doctor, LAMA = left against medical advice, or Absconded = disappeared during the process and are neither LWBS nor LAMA). Finally, another measure that has been gaining prominence is readmission rates. It is being used as a proxy for clinical quality of care, and we further discuss it in Section 6.1.1: out of the 2004–2005 ED patients, around 37% were eventually readmitted; and, overall, 3%, 11%, and 16% of the patients returned within 2, 14, and 30 days, respectively.

3.1.1. Research opportunities: Performance metrics. The ED performance metrics discussed above are mostly related to ED congestion. Hwang et al. (2011) lists over 70 such measures, which have given rise to prevalent crowding indices that support daily ED management (e.g. Bernstein et al. (2003); Hoot et al. (2007)). Such indices arose from ad-hoc statistical analysis that seeks to summarize (e.g. via regression) the state of ED congestion. Operations Research and Queueing Theory could complement these efforts by providing a natural environment for rigorously studying and developing congestion indices. For instance, operational regimes (Section 4.3) and state-space collapse in heavy-traffic (Huang, Carmeli and Mandelbaum (2015)) could yield rigorous state-summaries and sufficient statistics.

A practical challenge is that useful metrics are often difficult or impossible to measure from data. For example, our own data does not cover time-tillfirst-consultation, which is often part of triage protocols: e.g. following the

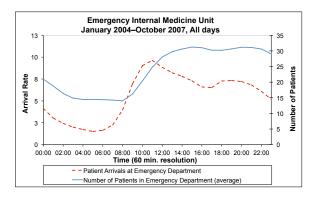


Fig 3. Average number of patients and arrival rate by hour of the day.

Canadian Triage and Acuity Scale (Canadadian-Triage), 90% of Category 3 (Urgent) patients should be seen by a physician within 30 minutes of arrival. The second example is patients' (im)patience (the *time* a patient is willing to wait before abandoning the ED), which is a natural building block for ED queueing models: while the overall abandonment proportion is observable, exact times till abandonment need not be. To be concrete, some patients notify the system about their abandonment; the others are either served, in which case their waiting time provides a lower bound for their patience, or they are discovered missing when called for service, which provides an upper bound. Statistical inference of ED (im)patience requires novel models and methods: these may draw from current-status (Sun, 2006) and survival analysis (Brown et al., 2005). Thus, estimating patients' (im)patience is an example of the research challenge to infer unobservable metrics from the measurable ones.

3.2. EDA: Time dependency and overdispersion. ED hourly arrival rates vary significantly over the day—see Figure 3, where it varies by a factor of over 5. We also observe a time-lag between the arrival rate and occupancy levels: it is due to the arrival rate changing significantly during a patient LOS, and it is formally explained by the time-varying version of Little's Law (Bertsimas and Mourtzinou, 1997). This lag, and in fact the daily shape of the arrival rate, clearly must be taken into account by staffing policies (Feldman et al., 2008; Green, Kolesar and Whitt, 2007; Yom-Tov and Mandelbaum, 2014).

Analyzing our data, Maman (2009) found support for the daily arrivals to fit a time-varying Poisson process, but with heterogeneity levels across days that render *random* the *arrival rate* itself. Kim and Whitt (2014) identified

similar patterns in a large Korean hospital. Such *overdispersion* (relative to Poisson arrivals) has significant operational consequences—the higher the variability the more pronounced is the effect (Koçağa, Armony and Ward (2015); Maman (2009)).

- 3.2.1. Research opportunities. Consider the time-varying shape of the arrival rate in Figure 3. Such temporal variability is typical of service systems, for example call centers (Brown et al., 2005), and it gives rise to two research questions: what are the main drivers of this shape, and how can a hospital affect it to benefit its patient flow?
- 1. What drives the shape of ED arrival rates? This question has not been systematically addressed. Its answer, which is a prerequisite for answering the second question above, could start with classifying shape drivers into natural, financial or behavioral. An example for a natural driver is that some time-periods are more emergency-prone than others; thus, there are relatively few arrivals during 2am-6am, and multi-trauma arrival rates (not depicted here) exhibit an early evening peak (and no morning peak)—conceivably at the time-of-day that is most vulnerable to such emergencies.

An (Israeli) example for a *financial* driver is that a referral letter from a Primary Care Physician (PCP) significantly reduces hospital charges for ED patients; consequently, most ED arrivals visit their PCP first, and since PCPs start seeing patients around 8am, these patients will start arriving to the ED around 10am. Interestingly, emergency maternity arrivals (again not depicted) peak at 9am; indeed, maternity patients do not need PCP referrals.

Behavioral drivers reflect preferences for some periods relative to others, which could add an explanation for the 10am peak: if patients are able to choose their time of travel, they would try to avoid the morning rush-hour.

2. Can a hospital affect the shape of its arrival-rates? The hospital has little control over arrivals due to natural factors, which hence must be accommodated by time-varying staffing (Green, Kolesar and Whitt, 2007; Yom-Tov and Mandelbaum, 2014) or ambulance diversion (Hagtvedt et al., 2009). At the same time, it may be able to affect the behavioral as well as the financial drivers of arrivals. Conceivably, both are associated with the less- or non-emergency patients, who enjoy (clinical) freedom in choosing their time of arrival. The shape of their arrival rate would hence favor convenience over need (in contrast to multi-trauma): this makes it a prime candidate for change, so that it fits the hospital's operational priorities. Indeed, hospitals have started making 'appointments' to ED visits (e.g. Gorman and Colliver (2014)), in an effort to balance workload.

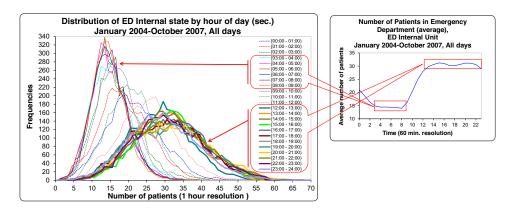


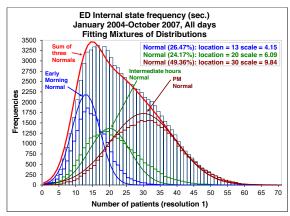
FIG 4. Internal ED Occupancy histogram (left) and Average Census (right), by hour of the day.

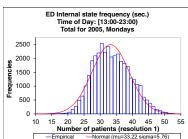
A related effort is control of temporary overloading via the provision of congestion information, for example predicted waiting time (Dong, Yom-Tov and Yom-Tov, 2014). The challenge here is to accurately forecast congestion, which has become an active area of research, in EDs (Plambeck et al., 2015) and elsewhere (Ibrahim and Whitt, 2011; Senderovich et al., 2015).

Indeed, many hospitals already advertise their ED waiting-time, and some do so by announcing a 4-hour moving average (e.g. HCA North Texas). This is not very informative as the ED state may change dramatically during 4-hour periods (see Figure 3). One could attempt to relate this uninformative granularity to the cheap-talk literature (e.g. Allon, Bassamboo and Gurvich (2011)) to help shape these announcements.

Finally, an indirect effort to affect workload shape is to identify, as early as possible, the least- or non-emergency patients, which is at the heart of any triage system. One could then possibly route the least-urgent to a fast-track and prepare the system, in advance, for those predicted to be hospitalized (Barak-Corren, Israelit and Reis, 2013).

3.3. EDA: Fitting a simple model to a complex reality. Figure 4 is details-heavy yet highly informative. Its left part shows 24 patient-count histograms for internal ED patients, each corresponding to a specific hour of the day, with reference (right) to mean patient count, also by hour of the day. (Similar shapes arise from total ED patient count—see Figure 10 in EV.) The figure displays a clear time-of-day behavior: There are two distinct bell-shaped distributions that correspond to low occupancy (15 patients on average) during the AM (3–9AM), and high (30 patients) during the PM (12–11PM); with two transitionary periods of low-to-high (9AM–12PM) and high-to-





- (a) Fitting a mixture of three Normal distributions (b) Fitting a Normal distribution to the Empirical distribution of ED occupancy
- for a specific year, day of the week, and time of day

FIG 5. Fitting parametric distributions to the Empirical distribution of ED occupancy.

low (11PM-3AM). We refer to these four periods as the four "occupancy regimes". Interestingly, when attempting to fit a mixture of three normal distributions to the ED occupancy distribution, SEEStat automatically detects the low, high and transitionary phases (See Figure 5a).

Further EDA (Figure 5b) reveals that, during peak times (PM), when controlling for factors such as day-of-the-week, patient type and calendar year, one obtains a good fit for the empirical distribution by a "steady-state" normal distribution with equal mean and variance. Hence, one might speculate that the underlying system dynamics can be modeled by an $M/M/\infty$ queue, which has a Poisson steady-state (mean = variance). Alternatively, however, it may also be modeled by an M/M/N+M queue with equal rates of service and abandonment (LWBS, LAMA, or Absconded). It follows that one cannot conclusively select a model through its empirical steady-state distribution (Whitt, 2012).

3.3.1. Research opportunities. In light of the above, one is led to seek the relevance-boundary of "black-box" ED models: they may support operational decisions that depend only on total patient count but not on internal dynamics (nor may these decisions alter internal dynamics); or they can model ED sojourn times within a larger hospital model. If in addition, and following Whitt (2012), a birth-death steady-state model is found appropriate for the "black-box" Dong and Whitt (2014), then model reversibility could accommodate applications that *change* total count: for example, am-

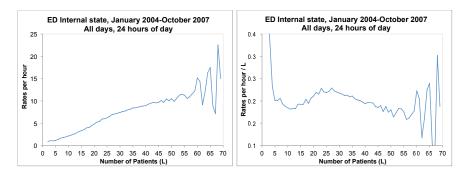


FIG 6. Service rate and service rate per patient as a function of number of patients.

bulance diversion when total count exceeds a certain threshold, which then truncates the count to this threshold (and the steady-state distribution is truncated correspondingly; see Chapter 1 in Kelly (1979)). On the other hand, such black-box models cannot support ED staffing (e.g. Yom-Tov and Mandelbaum (2014) acknowledges some internal network dynamics), or ambulance diversion that depends on the number of boarding patients.

In contrast to the macro level of our black-box model, one could consider a detailed model (such as simulation (Zeltyn et al., 2011)), which acknowledges explicitly micro-events at the level of individual patients and providers (physicians, nurses). The macro- and micro-models are two extreme cases of model granularity, with a range of levels in between (Huang, Carmeli and Mandelbaum, 2015; Marmor et al., 2012; Yom-Tov and Mandelbaum, 2014); The granularity level to be used depends on the target application, data availability and analytical techniques. Choosing the "right" level for an OR/queueing model has been mostly an art, which calls for systemizing this choice process. It could start with Whitt (2012) and Dong and Whitt (2014) that fit birth-death models, and continue with statistical techniques for model selection (e.g. Burnham and Anderson (2002)).

3.4. EDA: State dependency. In addition to time-dependent effects, we observe that the Internal ED displays some intriguing state-dependent behavior. Specifically, Figure 6 depicts service (or departure) rates as a function of the Internal patient count L: the graph on the left displays the total service rate, and the one on the right shows the service rate per patient. These graphs cannot arise from commonly-used (birth-death) queueing models such as M/M/N (for which the total departure rate is linearly increasing up to a certain point and then it is constant) or $M/M/\infty$ (for which a constant service rate per patient is expected). In contrast, the perpatient service rate has an interval $(11 \le L \le 20)$ where it is increasing

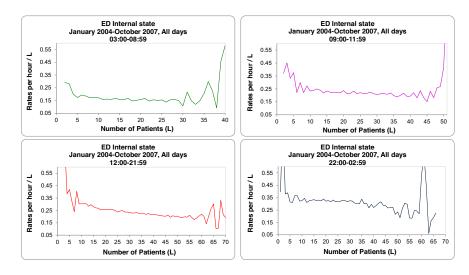


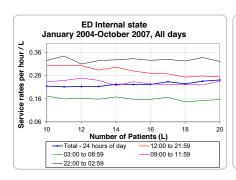
FIG 7. Service rate per patient as a function of L by occupancy regime.

in L, which is between two intervals of service-rate decrease. (The noise at the extremes, $L \leq 3$ and $L \geq 55$, is due to small sample sizes.) Note that Batt and Terwiesch (2014) and Kc and Terwiesch (2009) also found evidence for a state-dependent service rate.

Identifying the key factors that cause this particular state-dependence of the service rate per patient requires further exploration of the data. We start with explaining the apparent "speedup" effect $(10 \le L \le 25)$, followed by discussion of the slowdown effect in §3.4.1. As it turns out, this supposedly speedup is actually an artifact of biased sampling due to patient-heterogeneity and time-variability (Marmor et al., 2013). To see this, we further investigate the departure rate per patient, as a function of the patient count, at four different time-of-day intervals (corresponding roughly to the four occupancy regimes identified in Figure 4). For each of these, we observe, in Figure 7, either a constant service rate or a slowdown thereof, but no speedup.

Now, the rate-per-patient in Figure 6 is a weighted average of the four graphs of Figure 7. But these weights are not constant as a function of the patient count, as seen in Figure 8. Moreover, the service rate as a function of patient count varies at different times of the day. It follows that, what appears to be a speedup (increasing graph), is merely a weighted average of non-increasing graphs with state-dependent weights.

3.4.1. Research opportunities. As opposed to speedup, slowdown in service rate $(L \ge 25)$ appears to be real. Identifying the key factors that cause



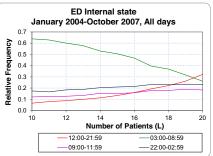


FIG 8. Service rate as a function of $10 \le L \le 20$ (left), and Relative frequency (weight) of occupancy regime per L (right).

this slowdown requires further research. We propose some plausible explanations next.

- Multiple resource types with limited capacity: As the number of occupied beds increases, the overall load on medical staff and equipment increases as well. Assuming a fixed processing capacity, the service rate per bed must then slow down.
- Psychological: Medical staff could become emotionally overwhelmed, to a point that exacerbates slowdown (Sullivan and Baghat, 1992).
- Choking: Service slowdown may also be attributed to so-called resource "choking": medical staff becomes increasingly occupied with caring for boarding ED patients (who create work while they wait and, moreover, their condition could actually deteriorate), which might end up taking capacity away from the to-be-released patients, thereby "choking" their throughput (see Figure 13 in Section 5.3). The choking phenomenon is well known in other environments such as transportation (Chen, Jia and Varaiya, 2001) and telecommunications (Gerla and Kleinrock, 1980), where it is also referred to as throughput degradation.
- Time dependency and patient heterogeneity: Finally, similar to the speedup effect, slowdown may also be attributed to the combination of time dependent arrivals and heterogenous patient mix.

In light of the above, one would like to identify the dominant factor that causes service-rate to slow down. Consequently, it is important to explore what can be done to alleviate this slowdown.

4. Internal wards. Internal Wards (IWs), often referred to as General Internal Wards or Internal Medicine Wards, are the "clinical heart" of a

 $\begin{array}{c} \text{Table 1} \\ \textit{Internal wards: operational profile} \end{array}$

	Ward A	Ward B	Ward C	Ward D	Ward E
Average LOS (days)	6.0	3.9	4.9	5.1	3.7
(STD)	(7.9)	(5.4)	(10.1)	(6.6)	(3.3)
Mean occupancy level	97.7%	94.4%	86.7%	96.9%	103.2%
Mean # patients per month	206.3	193.5	209.7	216.5	178.7
Standard (maximal)	45 (52)	30 (35)	44 (46)	42 (44)	24
capacity (# beds)					
Mean # patients per bed	4.58	6.45	4.77	5.16	7.44
per month					
Readmission rate	10.6%	11.2%	11.8%	9.0%	6.4%
(within 1 month)					

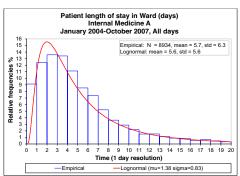
Data refer to period May 1, 2006–October 30, 2007 (excluding the months 1-3/2007, when Ward B was in charge of an additional 20-bed sub-ward).

hospital. Yet, relative to EDs, Operating Rooms and Intensive Care Units, IWs have received less attention in the Operations literature; this is hardly justified. IWs and other medical wards offer a rich environment in need of OR/OM research, which our EDA can only tap: It has revealed multiple time-scales of LOS, intriguing phenomena of scale diseconomies and coexisting operational-regimes of multiple resource types (beds, physicians). These characteristics are attributed to IW inflow design, capacity management and operational policies (e.g. discharge procedures, physician rounds).

4.1. Basic facts. Rambam hospital has five Internal Wards consisting of about 170 beds that cater to around 1000 patients per month. Wards A through D are identical from a clinical perspective; the patients treated in these wards share the same array of clinical conditions. Ward E is different in that it admits only patients of less severe conditions. Table 1 summarizes the operational profiles of the IWs. For example, bed capacity ranges from 24 to 45 beds and Average LOS (ALOS) from 3.7 to 6 days.

IWs B and E are by far the smallest (least number of beds) and the "fastest" (shortest ALOS, highest throughput). The short ALOS in IW E is to be expected as it treats the clinically simplest cases. In contrast, the "speed" of IW B is not as intuitive because this ward is assigned the same patient mix as IWs A,C, and D.

A shorter ALOS could reflect a more efficient clinical treatment or, alternatively, a less conservative discharge policy. Either must not arise from clinically premature discharges of patients, which would hurt patients clinical quality of care. To get a grasp on that quality, we use its operational (accessible hence common) proxy, namely patient readmission rate (proportion of



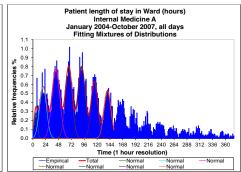


FIG 9. LOS distribution of IW A in two time-scales: daily and hourly.

patients who are re-hospitalized within a pre-specified period of time: one month in our case). In Table 1 we observe that the readmission rate of IW B is comparable to the other comparable wards (A–D). Moreover, patient surveys by Elkin and Rozenberg (2007) indicated that satisfaction levels do not differ significantly across wards. We conclude that IW B appears to be operationally superior yet clinically comparable to the other wards. This fact may be attributed to the smaller size of IW B, which we return to in Section 4.3.3.

4.2. EDA: LOS—A story of multiple time scales. Next, we examine the distribution of LOS in the IWs. While it is to be expected that clinical conditions affect patients LOS, the influence of operational and managerial protocols is less obvious. It turns out that some of this influence can be uncovered by examining the LOS distribution at the appropriate time scale.

Figure 9 shows the LOS distribution in IW A, in two time scales: days and hours. At a daily resolution, the Log-Normal distribution turns out to fit the data well. When considering an hourly resolution, however, a completely different distribution shape is observed: there are peaks that are periodically 24 hours apart, which correspond to a *mixture* of daily distributions. (We found that a normal mixture fits quite well, as depicted by the 7 normal mixture-components over the range of 0–150 hours in the right diagram of Figure 9.)

These two graphs reveal the impact of two operational protocols: The daily time scale represents physician decisions, made every morning, on whether to discharge a patient on that same day or to extend hospitalization by at least one more day. The second decision is the hour-of-day at which the patient is actually discharged. This latter decision is made according to the following discharge process: It starts with the physician who writes the dis-

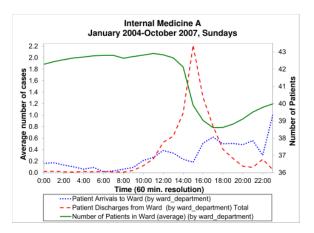


FIG 10. Arrivals, departures, and average number of patients in Internal wards by hour of day.

charge letters (after finishing the morning rounds); then nurses take care of paperwork, instructing patients on how to continue medical treatment after discharge, and then arranging for transportation (if needed). The discharge procedure is performed over "batches" of patients and, hence, takes a few hours. The result is a relatively low variance of the discharge time, as most patients are released between 3pm and 4pm—see Figure 10; which provides an explanation for the observed peaks in the hourly LOS distribution that are spaced 24 hours apart. The variation around these peaks is determined by the arrival process: patients are admitted into the IWs almost exclusively over a 12-hour period (10am–10pm) (Figure 10). Similar observations in a Singapore hospital led Shi et al. (2014) to model an inpatient ward as a 2-time-scale system, and to consequently propose flow-stabilization as a means of reducing delays.

We also observed two periods of unusual increase in arrival rate: one between 3pm-5pm, and a second towards midnight. The first is due to the phenomenon discussed above—patients discharge peaks at 3pm-4pm, and this enables higher transfer rates from the ED to the IWs. The influence of such transfers on ED congestion is discussed in Section 5.6. One plausible explanation for the second peak, that occurs towards midnight, lies within the ED shift schedule. ED physicians who work the night shift (which starts at midnight) are typically less experienced and in many cases are not authorized to approve patient hospitalization. As a result, towards the end of the evening shift, the more experienced physicians "clear off the table" in the ED, by making hospitalization decisions for many of the relevant patients; some of these end up being transferred to the IW later that night.

Who is the Server in an IW queueing model? Operational time-resolutions, specifically days/hours and hours/minutes for IWs, correspond to the time scale by which service durations are naturally measured which, in turn, identifies a corresponding notion of "a server". For example, IW LOS resolution in days corresponds to conceptualizing beds as servers, which is relevant in determining ward capacity. This is the setup in de Bruin et al. (2009) and Bekker and de Bruin (2010) who assume (hyper-) exponential LOS. (Log-Normal service durations are yet to be accommodated by queueing models.) Another IW resolution is hours, which is appropriate with servers being nurses, physicians or special IW equipment; in that case, the service times are measured in minutes or parts of an hour.

4.2.1. Research opportunities: Workload characterization, protocol mining, flow control, and why Log-Normal.

Offered load, or workload: The offered load is the skeleton around which capacity (staffing in the case of personnel) is dimensioned (Green, Kolesar and Whitt, 2007). Consider nurses as an example. Their time-varying offered load results from both routine and special care, and it varies during the day for at least two reasons (Mandelbaum, Momcilovic and Tseytlin, 2012): (a) routine care depends linearly on patient count, which varies over a day (Figure 10), and (b) admission and discharge of patients require additional work beyond routine, and it is more frequent during some hours than others (Figure 10). Combining both of these time variations, it is clear that staffing levels must (and actually do) vary during the day, hence the importance of observing and understanding the system in hourly resolution. As mentioned above, some efforts to develop queueing models for nurse staffing in medical wards have been carried out by Jennings and de Véricourt (2011), Green and Yankovic (2011) and Yom-Tov (2010). However, these works neither explain or incorporate the LOS distribution observed in our data, nor do they distinguish between routine, admission, and discharge workload. Even such a distinction might not be rich enough: indeed, the hospital environment calls for a broader view of workload, which we discuss in Section 5.5.4.

LOS and protocols: LOS or Delay distributions encapsulate important operational characteristics, and can hence be used to suggest, measure or track improvements. Consider, for example, the hourly effect of IW LOS (Figure 9), which is due to IW discharge protocols. It calls for an effort in the direction of smoothing IW discharge rates over the day (Shi et al., 2014). Taking an example from elsewhere at the hospital, consider the differences in shape of LOS distribution between two Maternity wards (§4.2.1 in EV), which result from differing patient mix; it suggests the redesign of routing

protocols towards a more balanced workload (Plonski et al., 2013). Queueing models are natural for analyzing the interplay between LOS distributions and operational protocols. This leads to open data-based questions in two directions: first, incorporating protocols (e.g. patient priorities, resource scheduling) in queueing models and validating the theoretical LOS distribution against data (performance); second and conversely, mining protocols from data. We now give two examples, one for each of the two directions.

Flow control: How will changes in the IW discharge process influence the system? For example, would the balancing of discharges, more uniformly over the day, benefit the entire hospital? How would such a change influence delays of patients waiting to be transferred into the IW from the ED? This connection between ED boarding and ward discharges was explored by Shi et al. (2014). We return to it in Section 5.6.

Why Log-Normal? A long-standing challenge is to explain the prevalence of Log-Normal as a distribution of service durations (e.g. IW LOS in days here, or durations of telephone calls in Brown et al. (2005)). Is Log-normality due to service protocols? Is it perhaps an inherent attribute of customer service requirements? Note that Log-Normal has an intrinsic structure that is both multiplicative—its logarithm is a central limit, and additive—it is infinitely divisible, being an integral against a Gamma process (Thorin, 1977). Can these properties help one explain the empirical Log-Normal service time distribution?

4.3. EDA: Operational regimes and diseconomies of scale. An asymptotic theory of many-server queues has developed and matured in recent years (Gans, Koole and Mandelbaum (2003) can serve as a starting point). This theory which has highlighted three main operational regimes: Efficiency Driven (ED), Quality Driven (QD) and Quality & Efficiency Driven (QED). The ED-regime prioritizes resource efficiency: servers are highly utilized (close to 100%), which results in long waits for service. In fact, waiting durations in the ED-regime are at least in the order of service times. In the QD-regime, the emphasis is on the operational quality of service: customers hardly wait for service, which requires that servers be amply staffed and thus available to serve. Finally, the QED-regime carefully balances service quality and server efficiency, thus aiming at high levels of both, and achieving it in systems that are large enough. Under the QED-regime, server utilization could exceed 90% while, at the same time, possibly half of the customers are served without delay, and those delayed wait one order of magnitude less than their service duration. Under some circumstances (e.g. Erlang-C

model), the QED-regime exhibits economies of scale in the sense that, as the system grows, operational performance improves.

Many-server queueing theory is based on asymptotic analysis, as the number of servers grows indefinitely. Nevertheless, QED theory has been found valuable also for small systems (few servers) that are not exceedingly overloaded. This robustness to system size is due to fast rates of convergence (Janssen, van Leeuwaarden and Zwart, 2011) and, significantly, it renders QED theory relevant to healthcare systems (Jennings and de Véricourt, 2011; Yom-Tov and Mandelbaum, 2014). One should mention that, prior to the era of many-server theory, asymptotic queueing theory was mostly concerned with relatively small systems—that is few servers that are too overloaded for QED to be applicable (e.g. hours waiting time for service times of minutes). This regime is nowadays referred to as conventional heavy-traffic (Chen and Yao, 2001) and, at our level of discussion, it will be convenient to incorporate it into the ED-regime.

In the following subsection, we seek to identify the operational regime that best fits the IWs. We then investigate (§4.3.3) how ward performance depends on its size. We shall argue that, although IW beds plausibly operate in the QED-regime, there is nevertheless evidence for *diseconomies* of scale.

4.3.1. In what regime do IWs operate? Can QED- and ED-regimes co-exist?. We start by identifying the operational regimes that are relevant to the IWs. These units have multiple types of servers (beds, nurses, physicians, medical equipment), that must be all considered. Here we focus on beds and physicians.

We argue that IW beds operate (as servers) in the QED-regime. To support this statement, we first note that our system of IWs has many (10's) beds/servers. Next we consider three of its performance measures: (a) bed occupancy levels; (b) fraction of patients that are hospitalized in non-IWs while still being under the medical care of IW physicians (patients who were blocked from being treated in IWs due to bed scarcity); (c) ratio between waiting time for a bed (server) and LOS (service time).

Considering data from the year 2008, we find that 3.54% of the ED patients were blocked, the occupancy level of IW beds was 93.1%, and patients waited hours (boarding) for service that lasted days (hospitalization). Such operational performance is QED—single digit blocking probability, 90+% utilization and waiting duration that is an order of magnitude less than service. Preliminary formal analysis, carried out in Section 4.3.1 of EV, demonstrates that QED performance of a loss model (Erlang-B, as in de Bruin et al. (2009)) usefully fits these operational performance measures of the IWs.

Turning to physicians as servers, we argue that they operate in the ED-regime (conventional heavy traffic). This is based on the following observation: from 4pm to 8am the following morning, there is a single physician on duty in each IW, and this physician admits the majority of new patients of the day. Therefore, patients that are admitted to an IW (only if there is an available bed) must wait until both a nurse and the physician become available. The admission process by the physician lasts approximately 30 minutes, and waiting time for physicians is plausibly hours (it takes an average of 3.2 hours to transfer a patient from the ED to the IWs; see Section 5.2). Performance of physicians is therefore Efficiency Driven.

4.3.2. Research opportunities. We identified two operational regimes, the QED- and ED-regime, that coexist within the ED+IW. What queueing models and operational regimes can valuably capture this reality? Note that such models must accommodate three time scales: minutes for physician treatment, hours for transfer delays, and days for hospitalization LOS. Some questions that naturally arise are the following: How do the regimes influence each other? Can we assume that the "bottleneck" of the system is the Efficiency Driven resource (physicians)? Thus, can one conclude that adding IW physicians is necessary for reducing transfer delays, while adding IW beds would have only a marginal impact on these delays? How would a change of IW physician priority influence the system, say giving higher priority to incoming patients (from the ED) over the already hospitalized (in the IWs)? Does the fact that the IW physicians operate in the ED-regime eliminate the economies of scale that one typically expects to find in QED systems? Empirical observations that will now be presented suggest that this might indeed be the case.

4.3.3. Diseconomies of scale (or how ward size affects LOS). Our data (Table 1) exhibits what appears to be a form of diseconomies of scale: a smaller ward (IW B) has a relative workload that is comparable to the larger wards, yet it enjoys a higher turnover rate per bed and a shorter ALOS, with no apparent negative influence on the quality of medical care. The phenomenon is reinforced by observing changes in LOS of IW B, when the number of beds in that ward changes. Figure 11 presents changes in ALOS and the average patient count, in IWs B and D over the years. During 2007, the ALOS of Ward B significantly increased. This was due to a temporary capacity increase, over a period of two months, during which IW B was made responsible for 20 additional beds. We observe that, although the same operational methods were used, they seem to work better in a smaller ward. In concert with the latter observation, we note a reduction in

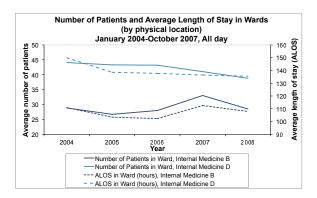


Fig 11. Average LOS and number of patients in Internal wards B and D by year.

ALOS of IW D, mainly from 2007 when ward size decreased as a result of a renovation. One is thus led to conjecture that there are some drawbacks in operating large medical units—e.g. larger wards are more challenging to manage, at least under existing conditions.

Indeed, several factors could limit the blessings of scale economies:

- Staffing policy: It is customary, in this hospital, to assign an IW nurse to a fixed set of beds; then nominate one experienced nurse to be a floater for solving emerging problems and help as needed. This setting gives little operational advantage to large units, if at all: the larger the unit the less a single floater can help each nurse. The tradeoff that is raised is between personal care (dedicated servers hence care continuity) vs. operational efficiency (pooling). This tradeoff has been addressed in call centers (Aksin, Karaesmen and Ormeci, 2007; Jouini, Dallery and Aksin, 2009), and in outpatient medical care (Balasubramanian, Muriel and Wang, 2012; Balasubramanian et al., 2010), but inpatient healthcare will surely add its own idiosyncracies. Another natural tradeoff that arises is whether the floater should indeed be an experienced nurse, or is it better to let more junior nurses be floaters so that they can learn from this broader experience.
- Centralized medical responsibility: Ward physicians share the responsibility over all patients. Every morning, the senior physicians, residents, interns, and medical students examine every patient case together (physician rounds) and discuss courses of treatment. This is essential as Rambam is a teaching hospital, and one of its central missions is the education and training of doctors. Naturally, the larger the unit the longer its morning round and, consequently, less capacity is available for other tasks (e.g. admissions and discharges).

4.3.4. Research opportunities. In Section 4.3.2 of EV, we provide additional plausible explanations for the observed diseconomies of scale. This phenomenon is important to model carefully and understand, as it can significantly affect decisions on unit sizing and operational strategy. While queueing theorists are well equipped to address the operational dimensions of such decisions, they must collaborate with researchers from other disciplines, such as organizational behavior, for a comprehensive treatment. This has been recently illustrated in Song, Tucker and Murrell (2015), that explores diseconomies of scale in the ED. They showed that when work is assigned to physicians early on, they develop an enhanced sense of ownership and tend to work faster, thereby negating positive effects of pooling. Such schemes of work assignment are discussed in §5.4 (see Figure 15).

Now, suppose one takes size differences among wards as a given fact (e.g. due to space constraints that cannot be relaxed). Then the following question arises: What protocol should be used to route patients from the ED to the wards, in order to fairly and efficiently distribute workload among them? This challenge is directly related to the process of transferring patients from the ED to the IWs, which is the topic of the next section.

- 5. The ED+IW network. After discussing the ED and IWs separately, in this section we discuss the ED+IW network as a whole. We start with the "ED-to-IW" process of transferring patients from the ED to the IWs: this is the "glue" that connects the ED to the IWs. We discuss delays in the transfer process (Sections 5.2–5.4) and fairness in this process towards both patients and medical staff (Section 5.5). We conclude, in Section 5.6, with an integrative view of the interplay between the three components: ED, IWs, and ED-to-IW.
- 5.1. ED-to-IW transfer process: Basic facts. The "ED-to-IW" process covers patient transfers from the ED to the IWs. We view this process in the context of flow or routing control. Routing in hospitals differs from routing in other service systems, for various reasons including incentive schemes, customers' (patients') limited control (or even helplessness), and the timing of the routing decision. Thus, although the transfer process involves routing-related issues similar to those that have been looked at extensively in the queueing literature, our data indicate that unusual system characteristics significantly affect delays and fairness features in a hospital setting, which creates many research opportunities.

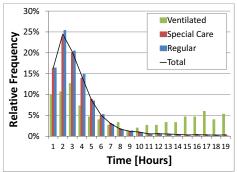
A patient, whom an ED physician decides to hospitalize in an IW, is assigned to one of five wards, according to a certain *routing policy* (described momentarily). If that ward is full, its staff may ask for reassignment with

the approval of the hospital's Head Nurse. Once the assigned ward is set, the ward staff prepares for this patient's arrival. In order for the transfer to commence, a bed and medical staff must be available, and the bed and equipment must be prepared for the specific patient (including potential rearrangement of current IW patients). Up to that point, the boarding patient waits in the ED and is under its care and responsibility. If none of the IWs is able to admit the patient within a reasonable time, the patient is "blocked", namely transferred to a non-internal ward. Then the latter undertakes nursing responsibilities while medical treatment is still provided by an IW physician.

An integral component of the transfer process is a routing policy, or patient assignment algorithm. As described in Section 4.2, Wards A–D provide similar medical services, while Ward E treats only the less severe patients. The similarity between Wards A–D requires a systematic assignment scheme of patients to these wards. Rambam hospital determines the assignment via a round-robin (cyclical) order among each patient type (ventilated, special care, and regular), while accounting for ward size (e.g. if Ward X has twice as many beds as Ward Y, then Ward X gets two assignments per one assignment of Y). This scheme is implemented by a computer software called "The Justice Table". As the name suggests, the algorithm was designed by the hospital to ensure fair distribution of patients among wards, so that staff workload will be balanced. It is worth noting that a survey among 5 additional hospitals in Israel (EV, Section 5.6) revealed that a cyclical routing policy is very common; still, some hospitals apply alternative assignment schemes, for example, random assignment by patient ID. Interestingly, only one of the surveyed hospitals uses an assignment that takes into account real-time bed occupancy.

5.2. Delays in transfer. As is customary elsewhere, the operational goal of Rambam hospital is to admit ED boarding patients to the IWs within four hours from decision of hospitalization. However, the delays are often significantly longer. The waiting-time histogram in Wards A–D, for the years 2006–2008, is depicted in Figure 12. We observe significant delays: while the average delay was 3.2 hours, 25% of the patients were delayed for more than 4 hours.

An interesting phenomenon is observed when analyzing transfer delays by patient type. We note that, on average, ventilated patients wait much longer (8.4 hours) than regular and special care patients (average of 3 and 3.3 hours, respectively)—see Figure 12. In particular, the delay distribution of these ventilated patients is bi-modal with 41% of such patients delayed



Patient	AVG	% delay	% delay
Type	(STD)	$\leq 4 \text{ h}$	> 10 h
Regular	3.00	77%	2%
	(2.53)		
Special	3.33	74%	5%
Care	(3.16)		
Ventilated	8.39	41%	41%
	(6.59)		
All Types	3.22	75%	4%
	(2.98)		

Data refer to period 5/06-10/08 (excluding the months 1-3/07, when Ward B was in charge of an additional sub-ward)

Fig 12. Transfer time by patient type, in hours.

by more than 10 hours. Ventilated patients must have the highest priority in transfer but, in reality, many do not benefit from it.

How come so many of the ventilated patients experience such long delays? We observe that the shorter delays of the ventilated patients (≤ 4 hours) have a pattern that resembles that of the other two patient types. The longer delays are harder to decipher. Possible explanations include: (a) Ventilated patients are hospitalized in a *sub-ward* inside the IW (A–D), often referred to as Transitional (intensive) Care Unit (TCU) (Armony, Chan and Zhu, 2013). Each such TCU has only 4–5 beds. The average occupancy rate of the TCUs at Rambam is 98.6%; the combination of high occupancy with a small number of beds results in long waits during overloaded periods. (b) Ventilated patients require a highly qualified staff to transfer them to the ward. Coordinating such transfers takes longer.

5.2.1. Research opportunities. Delays in transfer add opportunities to those arising from protocol mining, as discussed at the end of §4.2.1; relevant here is the specific challenge of deciphering a routing protocol from data such as in Figure 12. In addition, one would like to be able to analyze and optimize patient-flow protocols in queueing models, specifically here fork-join networks (representing synchronization between staff, beds and medical equipment) with heterogeneous customers. Such models, under the FCFS discipline, were approximated in Nguyen (1994). Their control was discussed in Atar, Mandelbaum and Zviran (2012) and Leite and Fragoso (2013).

The discussion above also raises again the tradeoff between the benefits of pooling and of continuity-of-care (see discussion of diseconomies of scale in §4.3.3). The fact that Rambam chose to distribute TCU beds among four IWs, instead of having one larger TCU, definitely increases waiting time for a TCU bed. Nevertheless, it is also advantageous from the quality-of-care perspective to have the TCU beds be part of an IW since, when patients' condition improves, they are transferred from the TCU in the IW to a regular room in the same IW, while continuing treatment by the same medical staff (physicians and nurses). This continuity-of-care reduces the number of hand-offs, which are prone to loss of information and medical errors. The tradeoff between pooling and continuity-of-care is an interesting challenge to navigate using OR methods.

5.3. Influence of transfer delays on the ED. Patients awaiting transfer (boarding patients) overload the ED: beds remain occupied while new patients continue to arrive, and the ED staff remains responsible for these boarding patients. Therefore, the ED in fact takes care of two types of patients: boarding patients (awaiting hospitalization) and in-process patients (under evaluation or treatment in the ED). Both types suffer from delays in the transfer process.

Boarding patients may experience significant discomfort while waiting: the ED is noisy, it is not private and often does not serve hot meals. In addition, ED patients do not enjoy the best professional medical treatment for their particular condition, and do not have dedicated attention as in the wards. Moreover, longer ED stays are associated with higher risk for hospital-acquired infections (nosocomial infections). Such delays may increase both hospital LOS and mortality rates, similarly to risks of delays in ICU transfer (e.g. Chalfin et al. (2007); Long and Mathews (2012); Maa (2011)). Hence, the longer patients wait in the ED, the higher the likelihood for clinical deterioration and the lower is their satisfaction.

In-process ED patients may suffer from delays in treatment, as additional workload imposed by boarding patients on ED staff can be significant. Figure 13 shows our estimates of the fraction of time that ED physicians spent caring for the boarding patients, assuming (consistently with the Rambam experience) that every such patient requires an average of 1.5 minutes of physician's time every 15 minutes. We observe that boarding patients take up to 11% of physician time in the ED. This extra workload for the ED staff, that occurs at times when their workload is already high, results in "wasted" capacity and throughput degradation, as discussed in Section 3.4.

To summarize, by improving patient flow from the ED to the IWs, in particular reducing boarding times, hospitals can improve the service and

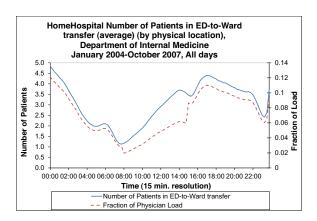


FIG 13. Number of patients in ED-to-IW transfer (A-E) and the fraction of time that ED physicians devote to these patients.

treatment provided to both transfer and in-process patients. In turn, reducing the workload in the ED would improve response to arriving patients and could, in fact, save lives.

- 5.3.1. Research opportunities. The delays in transfer give rise to the following research questions:
 - 1. Modeling transfer queue: Boarding patients may be viewed as customers waiting in queue to be served in the IW. Traditionally, in Queueing Theory, it has been assumed that customers receive service only once they reach a server, and not while waiting in queue. In contrast, here a waiting patient is "served" by both the ED and the IW. In the ED, clinical treatment is provided: according to regulations, boarding patients must be examined at least every 15 minutes. In the ward, "service" actually starts prior to the physical arrival of the patient, when the ward staff, once informed about a to-be-admitted patient, starts preparing for the arrival of this specific patient. The above has implications on modeling the ED-to-IW process, and it affects staffing, work scheduling, etc. A natural modeling framework here could be queueing networks with signals (Chao, Miyazawa and Pinedo, 1999).
 - 2. Emergency Department architecture: As described, ED staff attends to two types of patients: boarding and in-process. Each type has its own service requirements, leading to differing service distributions and differing distribution of time between successive treatments. While boarding patients receive periodic service according to a nearly-deterministic

schedule (unless complications arise), in-process schedule is random, in nature. Incorporating such a mix of distributions in queueing models and theory is a challenge of great interest.

The above distribution mix may also impact the choice of ED architecture; one may consider two options: (a) treating boarding and in-process patients together in the same physical location, as is done at Rambam, or (b) move the boarding patients to a transitional unit (sometimes called "delay room" or "observation room"), where they wait for transfer; this is done, for example, in a Singapore hospital that we were in contact with. Note that using option (b) implies having dedicated staff, equipment and space for this unit. The following question then arises: Under what conditions is each of these ED architectures more appropriate?

Interestingly, the Singapore architecture is even more complicated than (b) above, as the responsibility for the boarding patients is handed over to IW physicians after a two-hour ED boarding time. This provides the IW medical staff with an *incentive* to transfer the patients to the ward, as soon as possible, where they can be comfortably treated. In EV, Section 5.6, we further discuss how different architectures are related to incentive schemes and, in turn, influence delay times.

5.4. Causes of delay. In order to understand the causes of long delays in the ED-to-IW transfer, we interviewed hospital staff, conducted a time-andmotion study, and further explored our data. We learned that delays are caused not only by bed unavailability; patients often wait even when there are available beds. Indeed, our data shows that the fraction of patients who had an available bed in their designated ward, upon their assignment time, was 43%, 48%, 76%, 55%, for Wards A-D, respectively (which is consistent with our assertion in §4.3.1 that IW beds operate in the QED-regime). However, as Figure 12 shows, the probability to be admitted to the wards, immediately (or within a short time) after hospitalization decision, was much smaller. In fact, over the same period of time, only 4.9% of the patients were admitted to an IW within 30 minutes from their assignment to this ward. Our findings identify 13 plausible causes for delay, which are summarized in the Cause-and-Effect (Fishbone) diagram depicted in Figure 14. We elaborate here on two that have some interesting modeling aspects and present related research opportunities: 1) The timing of the routing decision, which is associated with the question of using Input-queued vs. Output-queued system, when routing patients from the ED to the IWs. And, 2) the influence on transfer delays of *information availability* during routing.

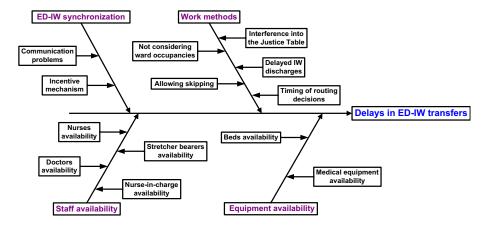


Fig 14. ED-to-IW delays—Cause and effect diagram.

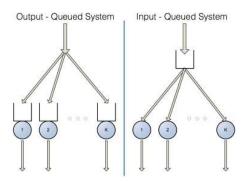


Fig 15. Output- vs. Input-queued scheme.

5.4.1. Research opportunities.

1. Timing of routing decision: Input-queued vs. Output-queued system. Recall that preparation for a transfer of a particular patient starts in the designated ward, prior to the actual transfer. This forces the hospital to adopt an output-queued scheme (Stolyar, 2005), where each patient is first assigned to an IW and then waits until the ward is able to admit. This is in contrast to a scheme where patients are placed in a "common" queue, then routed to an IW only once at the head of the line and a bed in any of the IWs becomes available. The latter is referred to as an input-queued scheme. Figure 15 depicts the two schemes.

Output-queued schemes are inherently less efficient than their inputqueued counterparts, because the routing decision is made at an earlier time with less information. Moreover, the output-queued scheme is inequitable towards patients because FCFS is often violated.

The problem of customer routing in input-queued schemes has received considerable attention in the queueing literature (e.g. Armony (2005); Atar and Shwartz (2008); Gurvich and Whitt (2010); Mandelbaum and Stolyar (2004)). Similar issues in output-queued systems have been generally overlooked. Exceptions include Stolyar (2005) and Tezcan (2008) who establish that the two systems have asymptotically similar performance, in both the conventional and the QED heavy-traffic regimes. This implies that inefficiencies, which arise in our ED-to-IW process due to the use of an output-queued scheme, become negligible in heavily-loaded systems. More generally, insights gained from studying the input-queued systems, as in the above references, may carry over to the output-queued systems. But how well does that insight translate to an environment such as a medical unit?

- 2. Not considering ward occupancies: The role of information availability in routing. An additional important aspect of routing schemes, which directly affects patient delays, is the availability of information on the system state, at the moment of the routing decision. On the one hand, hospitals may base the routing on no information, namely use a static routing policy like round robin (surprisingly, our experience suggests that this is a prevalent policy). On the other extreme, a full information policy that takes into account current occupancy levels and projected future dismissals and transfers is feasible, if the information system is accurate and accommodating enough (See Chapter 8 in Hall (2012)). Continuously-available accurate information is costly, however, and this cost-benefit tradeoff is yet to be explored.
- 5.5. Fairness in the ED-to-IW process. Transfer policies may have ramifications on fairness towards customers (patients) and towards servers (medical and nursing staff). We investigate both aspects next.
- 5.5.1. Fairness towards patients. In Section 5.4, we pointed out that output-queued schemes lead to diminished patient fairness, as FCFS order is often violated. (For references on the significance of FCFS in customer justice perception, see Mandelbaum, Momcilovic and Tseytlin (2012).) Indeed, our Rambam data indicate that 45% of the ED-to-IW transfer patients were "overtaken" by another patient (see Table 2). Moreover, more than a third of those were overtaken by at least three other patients. Although this figure includes overtaking between patient types, which may be due to clinical priorities, within each patient type there were significant FCFS violations

0 0		1 01		
IW \ Type	Regular	Special care	Ventilated	Total
Ward A	7.57%	7.33%	0.00%	7.37%
Ward B	3.86%	5.72%	0.00%	4.84%
Ward C	7.09%	6.62%	0.00%	6.80%
Ward D	8.18%	7.48%	2.70%	7.81%
Total within wards	6.91%	6.80%	0.67%	6.80%
Total in ED-to-IW	31%	31%	5%	

 $\begin{array}{c} {\rm TABLE} \ 2 \\ {\it Percentage} \ of \ FCFS \ violations \ per \ type \ within \ each \ IW \end{array}$

as well. Specifically, 31% were actually overtaken by at least one patient of their type, most of them not within the same ward; thus, these violations are conceivably due to the output-queued scheme.

While output-queues are inherently inefficient and unfair, they are unlikely to change in Rambam hospital due to the practical/clinical considerations described above, as well as psychological consideration (e.g., early ward assignment reduces uncertainty, which in turn reduces anxiety for patients and their families). The use of output-queues in the ED-to-IW process illustrates some idiosyncrasies of flow control in healthcare.

5.5.2. Research opportunities. A natural question is how to best maintain patient fairness in the output-queued scheme: What routing policies will keep the order close to FCFS? Is FCFS asymptotically maintained in heavy-traffic?

What other fairness criteria should be considered? Assuming that patients have preferences (clinical or prior experiences) for a specific ward, fairness may be defined with respect to the number of patients who are not assigned to their top priority. Related to this is the work of Thompson et al. (2009) that looks into minimizing the *cost* that reflects the number of "non-ideal" ward assignments; we propose to also look at the *equity* between patients in this context. One may alternatively consider achieving equity in terms of blocking probability (recall the discussion in §4.3.1) or patient delay.

5.5.3. Fairness towards staff. In Section 5.4, we discussed the implications of the routing policy on delays in the ED-to-IW process; in addition, routing also has a significant impact on wards' workload. High workload tends to cause personnel burnout, especially if work allocation is perceived as unjust (references can be found in Armony and Ward (2010)). Rambam hospital takes fairness into consideration, as is implied from the name "Justice Table". However, is the patient allocation to the wards indeed fair?

There are many candidates for defining server "fairness". One natural measure is equity in the occupancy level. Since the number of nurses and

doctors is typically proportional to the number of beds, equal occupancy levels imply that each nurse/doctor treats the same number of patients, on average. But does this imply that their workload is evenly distributed?

As mentioned in §4.2.1, staff workload in hospitals is not spread uniformly over a patient's stay, as patients admissions/discharges tend to be work intensive and treatment during the first days of a patient's hospitalization require much more time and effort from the staff than in the following days (Elkin and Rozenberg, 2007). Thus, one may consider an alternative fairness criterion: balancing the incoming load, or the "flux"—number of admitted patients per bed per time unit, among the wards. In Table 1 we observe that Ward B has a high average occupancy rate. In addition, as it is both the smallest and the "fastest" (shortest ALOS) ward, then (by Little's law) it has the highest flux among comparable IWs A–D. The workload of Ward B staff is hence the highest. We conclude that the most efficient ward is subject to the highest load—that is, patient allocation appears to be unfair towards ward staff.

Our data has already motivated some work on fair routing. Analytical results for input-queued systems were derived in Mandelbaum, Momcilovic and Tseytlin (2012), where both occupancy level and flux are taken into account with respect to fairness. Tseytlin and Zviran (2008) propose a simulation-supported algorithm that balances a weighted function of occupancy and flux to achieve both fairness and short delays in output-queued systems.

5.5.4. Research opportunities. In the context of output-queued systems, a more rigorous analytical study is needed to formalize the conclusions of Tseytlin and Zviran (2008). Specifically, how to combine the occupancy and flux criteria into a single effective workload measure, which would be balanced across wards. Even in the context of input-queued systems, it is our view that Armony and Ward (2010); Mandelbaum, Momcilovic and Tseytlin (2012) and Ward and Armony (2013) have just taken the first steps towards staff fairness, as they do not fully account for the dynamic nature of workload in healthcare. As patients progress in their hospital stay, their medical needs change (mostly reduce) and the accuracy in which one can predict their LOS increases. This information could be very useful in successfully balancing workload.

The underlying definition of operational fairness in our discussion thus far, proposed equal workload division across medical staff. A prerequisite for solving the "fairness problem" is then to define and calculate workload appropriately. However, we argue that such calculations must include not only direct time per resource but also emotional and cognitive efforts, as well as other relevant factors. For example, 1-minute of a standard chore does

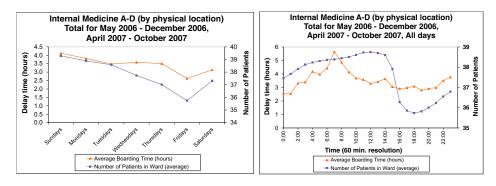


FIG 16. ED-to-IW boarding times and number of patients in IW.

not compare with a 1-minute life-saving challenge (Plonski et al., 2013). Thus, the mix of medical conditions and patient severity should also be included in workload calculation. For the latter, it is not straightforward to determine whether wards would be inclined to admit the less severe patients (who add less workload, and potentially less emotional stress), as opposed to the more severe patients, who would challenge the medical staff, thus providing them with further learning and research opportunities; the latter is especially relevant in teaching hospitals such as Rambam.

5.6. A system view. In this Section we underscore the importance of looking at this network of ED, IWs and ED-to-IWs as a whole, as these three components are clearly interdependent. For concreteness, we focus on how the discharge policy in the IW affects ED-to-IW transfer times which, in turn, affect ED workload. We thereby argue that an integrative system view is appropriate.

It is natural to expect that the higher the occupancy in the IWs the longer the boarding times, due to limited IW resources. The left diagram in Figure 16 displays the average boarding time alongside the average number of patients per ward—in IWs A–D, by day of the week. We observe that, as expected, the two measures have a similar weekly pattern. The right diagram in Figure 16 shows delays in the transfer process and the average number of patients in the IWs, as they vary throughout the day. The correlation here is not as apparent as in the daily resolution; other factors, such as the IW discharge process, also play a role.

We observe that the longest delays are experienced by patients assigned to the IWs in early morning (6am–8am)—these patients need to wait on average 5 hours or more. This is due to a combination of two reasons. First, as the night progresses, most of the IW beds become occupied (see Figure 10).

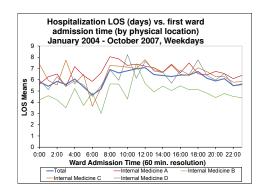


FIG 17. ALOS in IWs A through D, as a function of ward admission-time.

Hence, many of the patients arriving in the morning are unlikely to encounter an available bed, while new beds will only become available towards 3pm. A second reason stems from the working procedures of physicians within the IW. As discussed in Section 5.4 (see Figure 14), an available physician is required to complete a transfer. But most of the IW physicians perform their morning rounds at this time and cannot admit new patients. We note a consistent decline in the transfer delay up until noon. Patients assigned to the IWs during these times are admitted into the IWs between 1pm-3pm. This is about the time when the physicians' morning rounds are complete; staff and beds are starting to become available. Indeed, there is a sharp decline in the number of IW patients around 3pm-4pm when most of the IWs discharges are complete.

Further data analysis reveals that patients who are admitted to the IWs before 8am experience a significantly shorter LOS; Figure 17 shows that early hospitalization may reduce ALOS by more than 1 day. A correlation between hospitalization time and LOS was also reported by Earnest, Chen and Seow (2006): they observed that patients who are admitted in afternoon/night hours have ALOS that is longer than patients admitted in the morning. In contrast, we differentiate between early- and late-morning admissions. Regardless, in both cases, the plausible explanation for the difference in ALOS is the same: If patients are admitted to the ward early enough, the first day of treatment is more effective, as tests, medical procedures and treatments start earlier, and hence LOS is reduced. Thus, we argue that it is important to shorten the ED-to-IW transfer process and improve the IW admission process so that the first day of hospitalization is not "wasted".

In Section 5.3, we discussed how boarding times impact physician workload in the ED and hence may influence quality of care there. Thus, we

observe a chain of events in which the discharge policy in the IWs impacts the delays in transfer, which in turn affects workload in the ED. In particular, a system-view perspective is called for.

5.6.1. Research opportunities. Our discussion suggests that daily routines (schedules) in the IWs have a significant impact on boarding times and thereby on ED workload. At the same time, these routines also affect IW LOS. The question arises as to if and how one might wish to change daily routines in view of these effects. The question fits well within a queueing context. The present daily routine at Rambam may be viewed as a priority scheme where currently-hospitalized IW patients enjoy priority during morning physicians' rounds; these patients become low-priority, as discharged patients obtain priority in the afternoon, then followed by newly-admitted patients. Is it possible to positively affect overall system performance by altering these priorities (e.g. prioritizing newly-admitted and to-be discharged patient in the morning)? More broadly, the challenge falls within the uncharted territory of designing priority schemes for time-varying queueing networks.

Our discussion here brings us back to the broader issue—that is the need for a system view, in order to understand and relieve delays in patient flow. Consider, for example, the boarding patients in EDs (Figure 16) or in ICUs (Long and Mathews, 2012). Long boarding times are often due to scarce resources or synchronization gaps (Zaied, 2011), which are rooted in parts of the system that differ from those where the delays are manifested. For example, scarce resources in the IWs exacerbate ED delays, and tardy processing of MRI results can prolong ICU LOS. It follows that a system view is required for the analysis of patient flow in hospitals.

When analyzing ED+IWs flows (§5), the wards operate naturally on a time-scale of days while the ED time scale is hours. Wards thus serve as a random environment for the ED (Ramakrishnan, Sier and Taylor, 2005). Figure 9 (§4.2) reveals that the hourly scale is also of interest for IWs. These empirical observations arise in a service system (hospital) that evolves in multiple time scales, which are all natural and relevant for measuring and modeling its performance. The mathematical manifestation of such scales is asymptotic analysis that highlights what matters at each scale, while averaging out details that are deemed insignificant (e.g., Mandelbaum, Momcilovic and Tseytlin (2012), Shi et al. (2014), Gurvich and Perry (2012) and Zacharias and Armony (2013)).

6. Discussion and concluding remarks. We have described research opportunities that arose from EDA of operational patient flow data. We

now discuss the relationship between operational performance measures and overall hospital performance, and conclude with comments on data-based OR research.

6.1. Operational measures as surrogates to overall hospital performance. Hospital performance is measured across a variety of dimensions: clinical, financial, operational, psychological (patient satisfaction) and societal. The most important measures are clearly clinical but, practically, operational performance is the easiest to quantify, measure, track and react upon in real time. Moreover, operational performance is tightly coupled with the other dimensions (e.g. rate of readmissions with quality of clinical care, or LOS and LWBS with financial performance), which explains its choice as a "language" that captures overall hospital performance.

Operational performance measures are often associated with patient flow. Among these, we discussed LWBS (Section 3) and "blocking" (where patients end up being hospitalized in a ward different from that which is medically best for them—Section 4.3.1); boarding (transfer) time from the ED to the appropriate medical unit; and measures related to LOS, in the ED or IWs, such as merely averages (or medians), or fractions staying beyond a desired threshold. Other measures that have not been mentioned here require intra-ward data, which is beyond our data granularity. Examples include the time until triage or until a patient is first seen by an ED physician (Zeltyn et al., 2011), the number of visits to a physician during an ED sojourn (Huang, Carmeli and Mandelbaum, 2015) and the time-ingredients of an ED visit (treatment and waiting—for a resource, for synchronization or for a treatment to take its effect; see Zaied (2011) and Atar, Mandelbaum and Zviran (2012)).

6.1.1. Readmissions. As already indicated in Sections 3.1 and 4.1, our data supports the analysis of readmissions (Mandelbaum et al., 2013). We now elaborate on this operational measure of performance since policy makers are increasingly focusing on it, as part of efforts to extend quality of care measures from within-hospital processes to after-hospital short-term outcomes (Medicare USA, 2013). As mentioned, the likelihood of readmission to the hospital, within a relatively short time, is a natural indirect measure for quality of care (similarly to first-call-resolution rates in call centers). Consequently, readmission rates are accounted for when profiling hospitals' quality and determining reimbursements for their services.

One should argue that readmissions should be considered judiciously as some of them could be due to factors outside hospital control, or they may be an integral part of the treatment regimen. For example, returns within a few months to chemotherapy are typically planned and are unrelated to poor quality. But there are also unplanned chemotherapy returns after 1–2 weeks, which arise from complications after treatment. To properly incorporate readmissions in a queueing model (such as in Yom-Tov and Mandelbaum (2014)) one should distinguish between these two readmission types by, for example, modeling planned (unplanned) readmissions as deterministic (stochastic) returns. Also note that readmissions should be measured in their natural time-scale. For example, readmission to an ED should be measured in a time scale of days-weeks, while readmissions to an IW have a natural time-scale of weeks-months.

6.1.2. Capacity and cost of care. Of utter importance to hospital managers and policy makers is hospital costing. Kaplan and Porter (2011) argue that the mapping of patient / process flow, and the association of its activities with their supporting resources, should constitute the first step in understanding the cost of care. This is consistent with promoting a queueing-network view for understanding and calculating cost of care. Indeed, Kaplan and Porter (2011) further submit that most hospital costs are mistakenly judged as fixed while they ought to be viewed as variable costs; it follows that the corresponding resource levels are in fact flexible, an observation that renders controllable most resources in a hospital.

This viewpoint naturally connects with the distinction between static and dynamic capacity in queueing systems, which we now explain. Capacity of a hospital or a ward is commonly expressed in terms of the number of beds (or rooms, or physical space). However, it is also necessary to associate with a ward its processing capacity, which is determined by its human and equipment resources: nurses, physicians, support personnel, and medical apparatus. One thus distinguishes between static capacity (e.g. beds) and dynamic (processing) capacity of a resource. (Note that bed capacity plays the dual role of static capacity—capping the number of patients that can be simultaneously hospitalized, and dynamic capacity—serving as a proxy for the processing capacity of medical personnel). This distinction connects back to costs in the sense that static capacity is thought of as fixed over the relevant horizon, hence its cost is fixed; processing capacity, on the other hand, is considered variable in that its level (and hence also cost) is flexible (controllable).

6.2. Some concluding comments on data-based research—A great opportunity but no less of a challenge. The goal of the present work has been two-fold: first, to encourage and strengthen, through data and its EDA, the natural link between queueing theory and its application to patient flow in

healthcare; and second, to facilitate data-based learning for researchers who seek to reinforce this important link.

While theory has been the comfort zone of Operations Research (OR) and Applied Probability (AP), the situation dramatically differs when data is brought into the picture. Fundamental changes are therefore essential—both within our OR/AP community as well as our potential healthcare partners: changes in accessibility to healthcare data, in education (e.g. concerning the necessity of data-based OR research, importance and need to publish EDA, benefits of research reproducibility) and in funding priorities (e.g. for developing and sustaining the infrastructure that is a prerequisite for a research such as the one reported here).

However, we are cautiously optimistic. Indeed, comprehensive data collection is becoming increasingly feasible, systematic and cheaper, for example via Real-time Location Systems (RTLS), which will ultimately integrate with Personal-Health and Financial Records. This will enable partnerships between researchers and healthcare providers—partnerships that are based on multidisciplinary (clinical, operational, financial, psychological) tracking, of complete care-paths that start possibly at onset of symptoms (rather than hospital entrance), and at resolution levels of the individual patient and provider. The process of data-based OR research in hospitals is thus only beginning¹.

Acknowledgements. We dedicate our paper to the memory of the late David Sinreich. As a professor at the Technion IE&M Faculty, David was a pioneer and an example-to-follow, in recognizing the necessity and importance of data- and simulation-based research in healthcare.

Our research, and its seed financial backing, started within the OCR (Open Collaborative Research) Project, funded by IBM and led jointly by IBM Haifa Research (headed by Oded Cohn), Rambam Hospital (CEO Rafi Beyar) and Technion IE&M Faculty (past dean Boaz Golany). Indeed, it was Beyar who invited us at the outset to "use" the healthcare campus as our research laboratory; this invitation has given rise to the present work and many offsprings.

Data analysis, maintenance of data repositories, and continuous advice and support have been provided by the Technion SEE Laboratory and its affiliates: Valery Trofimov, Igor Gavako, Ella Nadjharov, Shimrit Maman, Katya Kutsy, Nitzan Carmeli, and Arik Senderovic.

¹Complete operational coverage, within the hospital, is already feasible; SEELab is now continuously getting RTLS data from a large day-hospital, which is collected through 900 sensors spread over 7 floors: it covers around 900 patients and 300 providers daily, at 3-second resolution, and it amounts to more than 1GB of raw data per week.

Financially, the research of YM, GY and YT was supported by graduate fellowships from Technion's Graduate School and the Israel National Institute for Health Services and Health Policy Research. The joint research of MA and AM was funded by BSF (Bi-national Science Foundation) Grants 2005175/2008480. AM was partially supported by ISF Grant 1357/08 and by the Technion funds for promotion of research and sponsored research. MA was partially funded by the Lady Davis Fellowship as a visitor at the Technion IE&M faculty.

Some of AM's research was funded by and carried out while visiting SAMSI NSF, STOR UNC, IOMS Stern NYU and Statistics Wharton UPenn—the hospitality of all these institutions is truly appreciated.

Our paper greatly benefited from feedback by colleagues and a thorough and insightful refereeing process. Ward Whitt discovered a flaw in our empirical analysis that led to the rewriting of the ED section. Jim Dai read carefully the first version and provided significant and helpful editorial feedback. We are very grateful to Peter Glynn, founding editor-in-chief of *Stochastic Systems*, for his encouragement and support throughout. He was joined by an anonymous AE, to lead and guide us safely through the revision process.

Last but certainly not least, a project such as this one requires a collaborative effort by many individuals. Their list, which is too long to acknowledge here, must have at its top our host Rambam hospital: its management, nursing staff and physicians made as good a partner as one could hope for.

APPENDIX: A MODEL FOR OR/AP DATA-BASED RESEARCH

The traditional still prevalent model for data-based OR/AP research has been one where an *individual* researcher, or a small group, obtains and analyzes data for the sake of an *isolated* research project. Our experience is that such a model cannot address today's empirical needs. For example, hospital data is typically large, complex, contaminated and incomplete, which calls for a professional and inevitably time-consuming treatment. Next, using data in a single project, or a few for that matter, is wasteful. This calls for data-reuse and sharing, across student generations or research groups, which requires infrastructure, documentation, maintenance and coordination. The challenge is exacerbated by healthcare data being confidential and proprietary, which prevents reproducibility and slows down progress (Nestler, 2011).

A feasible model. An alternative model is a research laboratory, funded by and serving a small community of researchers, with access to its data being as broad as possible. A working example is the Technion SEELab, where

readers can access RambamData. Little effort will be then required to reproduce our present EDA and going beyond it. In fact, most of our figures were created by SEEStat—a SEELab-developed user-friendly platform for online (real-time) EDA—and readers can recreate this process by following Nadjhahrov et al. (2013).

SEELAB

To elaborate some, SEELab is a data-based research laboratory founded in 2007 and residing at the IE&M Faculty of the Technion in Haifa, Israel. (SEE stands for "Service Enterprise Engineering".) SEELab maintains a repository of transaction-level operational data (log-files) from large service operations. This data is collected and cleaned, thus preparing it for research and teaching. Currently, SEELab databases include transaction-level multi-year data from 4 call centers, an internet academic website, 8 emergency departments (mainly arrivals data), a large ambulatory hospital (RTLS-based), and 4 years of data from the Rambam Hospital—the latter is the empirical foundation for the present research.

The EDA environment of SEELab is SEEStat—a software platform that enables real-time statistical analysis of service data at seconds-to-months time resolutions. SEEStat was used to create most of our figures. It implements many statistical algorithms: parametric distribution fitting and selection, fitting of distribution mixtures, survival analysis and more—with all algorithms interacting seamlessly with all the databases. SEEStat also interacts with SEEGraph, a pilot-environment for structure-mining, on-demand creation, display and animation of data-based process maps (e.g. Figure 1, and the animation of its underlying data (SEEnimations)).

WORKING ONLINE WITH SEESTAT

Three SEELab databases are *publicly accessible* at SEEServer, the SEELab server: two from call centers and one from the Rambam hospital. The Rambam data is described in §1.2, and our analysis of it greatly benefitted from our call-center experiences.

The connection protocol to SEEStat, for any research or teaching purpose, is simply as follows: go to the SEELab webpage http://ie.technion.ac.il/Labs/Serveng; then proceed, either via the link SEEStat Online, or directly through http://seeserver.iem.technion.ac.il/see-terminal, and complete the registration procedure. Within a day or so, you will receive a confirmation of your registration, plus a password that allows you access to SEEStat, SEELab's EDA environment, and via SEEStat to the above-mentioned databases. Note that your confirmation email includes two

attachments: a trouble-shooting document and a self-taught tutorial that is based on call center data and the Rambam hospital data. We propose that you print out the tutorial, connect to SEEStat and then let the tutorial guide you, hands-on, through SEEStat basics—this should take no more than 1.5 hours.

Reproducing our EDA and beyond. Rambam data is publicly available, either for downloading (RambamData) or through SEEStat, as described above. The download link includes data documentation. To facilitate reproducibility, the document Nadjhahrov et al. (2013) provides a detailed description of the creation process of our EDA, which includes all figures (except for Figure 12) in the present paper.

REFERENCES

- AKSIN, O. Z., KARAESMEN, F. and ORMECI, E. L. (2007). A Review of Workforce Cross-Training in Call Centers from an Operations Management Perspective. In *Workforce Cross Training Handbook* (D. Nembhard, ed.), CRC Press.
- ALLON, G., BASSAMBOO, A. and GURVICH, I. (2011). "We Will Be Right with You": Managing Customer Expectations with Vague Promises and Cheap Talk. Operations Research 59 1382–1394. MR2872007
- Armony, M. (2005). Dynamic Routing in Large-Scale Service Systems with Heterogeneous Servers. *Queueing Systems* **51** 287–329. MR2189596
- Armony, M., Chan, C. W. and Zhu, B. (2013). Critical Care in Hospitals: When to Introduce a Step Down Unit? Working paper, Columbia University.
- Armony, M. and Ward, A. (2010). Fair Dynamic Routing in Large-Scale Heterogeneous-Server Systems. *Operations Research* **58** 624–637. MR2680568
- ARMONY, M., ISRAELIT, S., MANDELBAUM, A., MARMOR, Y. N., TSEYTLIN, Y. and YOM-TOV, G. B. (2015). On Patient Flow in Hospitals: A Data-Based Queueing-Science Perspective. An Extended Version (EV). Working paper, http://ie.technion.ac.il/serveng/References/Patient%20flow%20main.pdf.
- Atar, R., Mandelbaum, A. and Zviran, A. (2012). Control of Fork-Join Networks in Heavy Traffic. Allerton Conference.
- Atar, R. and Shwartz, A. (2008). Efficient Routing in Heavy Traffic under Partial Sampling of Service Times. *Mathematics of Operations Research* **33** 899–909. MR2464649
- Balasubramanian, H., Muriel, A. and Wang, L. (2012). The Impact of Flexibility and Capacity Allocation on the Performance of Primary Care Practices. *Flexible Services and Manufacturing Journal* **24** 422–447.
- Balasubramanian, H., Banerjee, R., Denton, B., Naessens, J., Wood, D. and Stahl, J. (2010). Improving Clinical Access and Continuity Using Physician Panel RSSedesign. *Journal of General Internal Medicine* **25** 1109–1115.
- BARAK-CORREN, Y., ISRAELIT, S. and REIS, B. Y. (2013). Progressive Prediction of Hospitalization in The Emergency Department: Uncovering Hidden Patterns to Improve Patient Flow. Working paper.
- BARON, O., BERMAN, O., KRASS, D. and WANG, J. (2014). Using Strategic Idleness to Improve Customer Service Experience in Service Networks. Operations Research 62 123–140. MR3188591

- BATT, R. J. and TERWIESCH, C. (2014). Doctors Under Load: An Empirical Study of State Dependent Service Times in Emergency Care. Working paper.
- BEKKER, R. and DE BRUIN, A. M. (2010). Time-Dependent Analysis for Refused Admissions in Clinical Wards. *Annals of Operations Research* **178** 45–65. MR2659096
- BERNSTEIN, S. L., VERGHESE, V., LEUNG, W., LUNNEY, A. T. and PEREZ, I. (2003). Development and Validation of a New Index to Measure Emergency Department Crowding. Academic Emergency Medicine 10 938–942.
- Bertsimas, D. and Mourtzinou, G. (1997). Transient Laws of Non-stationary Queueing Systems and Their Applications. *Queueing Systems* **25** 115–155. MR1458588
- Brandeau, M. L., Sainfort, F. and Pierskalla, W. P., eds. (2004). Operations Research and Health Care: A Handbook of Methods and Applications. Kluwer Academic Publishers, London.
- Brown, L., Gans, N., Mandelbaum, A., Sakov, A., Shen, H., Zeltyn, S. and Zhao, L. (2005). Statistical Analysis of a Telephone Call Center: A Queueing-Science Perspective. *Journal of the American Statistical Association* **100** 36–50. MR2166068
- Burnham, K. P. and Anderson, D. R. (2002). Model Selection and Multimodal Inference: A Practical Information-Theoretic Approach, 2nd Edition. Springer. MR1919620
- CANADADIAN-TRIAGE Admission of Paitents to Over-Capacity Inpatient Beds. Appendix A, http://www.calgaryhealthregion.ca/policy/docs/1451/Admission_over-capacity_AppendixA.pdf.
- Chalfin, D. B., Trzeciak, S., Likourezos, A., Baumann, B. M. and Dellinger, R. P. (2007). Impact of Delayed Transfer of Critically Ill Patients from the Emergency Department to the Intensive Care Unit. *Critical Care Medicine* **35** 1477–1483.
- Chan, C., Farias, V. and Escobar, G. (2014). The Impact of Delays on Service Times in the Intensive Care Unit. Working paper.
- Chan, C., Yom-Tov, G. B. and Escobar, G. (2014). When to Use Speedup: An Examination of Service Systems with Returns. *Operations Research* **62** 462–482. MR3209183
- Chao, X., Miyazawa, M. and Pinedo, M. (1999). Queueing Networks: Customers, Signals and Product Form Solutions. Wiley.
- CHEN, C., JIA, Z. and VARAIYA, P. (2001). Causes and Cures of Highway Congestion. Control Systems, IEEE 21 26–33.
- CHEN, H. and YAO, D. D. (2001). Fundamentals of Queuing Networks: Performance, Asymptotics, and Optimization. Springer. MR1835969
- COOPER, A. B., LITVAK, E., LONG, M. C. and McManus, M. L. (2001). Emergency Department Diversion: Causes and Solutions. *Academic Emergency Medicine* 8 1108–1110.
- DE BRUIN, A. M., VAN ROSSUM, A. C., VISSER, M. C. and KOOLE, G. M. (2007). Modeling the Emergency Cardiac In-Patient Flow: An Application of Queuing Theory. *Health Care Management Science* **10** 125–137.
- DE BRUIN, A. M., BEKKER, R., VAN ZANTEN, L. and KOOLE, G. M. (2009). Dimensioning Hospital Wards Using the Erlang Loss Model. *Annals of Operations Research* **178** 23–43. MR2659095
- Denton, B. T., ed. (2013). Handbook of Healthcare Operations Management: Methods and Applications. Springer.
- Dong, J. and Whitt, W. (2014). On Fitted Birth-and-Death Queue Models. Working paper, Columbia University.
- Dong, J., Yom-Tov, E. and Yom-Tov, G. B. (2014). Hospital Network Synchronization Through Waiting Time Announcements. Working paper.

- EARNEST, A., CHEN, M. and SEOW, E. (2006). Exploring if Day and Time of Admission is Associated with Average Length of Stay Among Inpatients from a Tertiary Hospital in Singapore: An Analytic Study Based on Routine Admission Data. *BMC Health Services Research* 6 6.
- ELKIN, K. and ROZENBERG, N. (2007). Patients Flow from the Emergency Department to the Internal Wards. IE&M project, Technion (In Hebrew).
- Feldman, Z., Mandelbaum, A., Massey, W. A. and Whitt, W. (2008). Staffing of Time-Varying Queues to Achieve Time-Stable Performance. *Management Science* **54** 324–338.
- FROEHLE, C. M. and MAGAZINE, M. J. (2013). Improving Scheduling and Flow in Complex Outpatient Clinics. In *Handbook of Healthcare Operations Management: Methods and Applications* (B. T. Denton, ed.) 9, 229–307. Springer.
- GANS, N., KOOLE, G. and MANDELBAUM, A. (2003). Telephone Call Centers: Tutorial, Review and Research Prospects. *Manufactoring, Services and Operations Management* 5 79–141.
- GERLA, M. and KLEINROCK, L. (1980). Flow Control: A Comparative Survey. *IEEE Transactions on Communications* **28** 553–574.
- GORMAN, A. and COLLIVER, V. (2014). The Latest In Medical Convenience: ER Appointments. Chronicle for Kaiser Health News. http://kaiserhealthnews.org/news/the-latest-in-medical-convenience-er-appointments/.
- GREEN, L. (2004). Capacity Planning and Management in Hospitals. In Operations Research and Health Care: A Handbook of Methods and Applications (M. L. Brandeau, F. Sainfort and W. P. Pierskalla, eds.) 14–41. Kluwer Academic Publishers, London
- GREEN, L. V. (2008). Using Operations Research to Reduce Delays for Healthcare. In *Tutorials in Operations Research* (Z.-L. Chen and S. Raghavan, eds.) 1–16. INFORMS.
- Green, L. V., Kolesar, P. J. and Whitt, W. (2007). Coping with Time-Varying Demand When Setting Staffing Requirements for a Service System. *Production and Operations Management* 16 13–39.
- GREEN, L. and YANKOVIC, N. (2011). Identifying Good Nursing Levels: A Queuing Approach. Operations Research 59 942–955. MR2844415
- Green, L., Soares, J., Giglio, J. F. and Green, R. A. (2006). Using Queuing Theory to Increase the Effectiveness of Emergency Department Provider Staffing. *Academic Emergency Medicine* **13** 61–68.
- Gurvich, I. and Perry, O. (2012). Overflow Networks: Approximations and Implications to Call-Center Outsourcing. *Operations Research* **60** 996–1009. MR2979436
- Gurvich, I. and Whitt, W. (2010). Service-Level Differentiation in Many-Server Service Systems via Queue-Ratio Routing. *Operations Research* **58** 316–328. MR2674799
- HAGTVEDT, R., FERGUSON, M., GRIFFIN, P., JONES, G. T. and KESKINOCAK, P. (2009). Cooperative Strategies To Reduce Ambulance Diversion. *Proceedings of the 2009 Winter Simulation Conference* 266 1085–1090.
- HALL, R. W., ed. (2012). Handbook of Healthcare System Scheduling. Springer.
- HALL, R. W., ed. (2013). Patient Flow: Reducing Delay in Healthcare Delivery. Springer. 2nd edition.
- HALL, R., BELSON, D., MURALI, P. and DESSOUKY, M. (2006). Modeling Patient Flows Through the Healthcare System. In *Patient Flow: Reducing Delay in Healthcare Deliv*ery (R. W. Hall, ed.) 1, 1–45. Springer.
- HOOT, N. R., ZHOU, C., JONES, I. and ARONSKY, D. (2007). Measuring and Forecasting Emergency Department Crowding in Real Time. *Annals of Emergency Medicine* **49** 747–755.

- Huang, J. (2013). Patient Flow Management in Emergency Departments. PhD thesis, National University of Singapore (NUS).
- HUANG, J., CARMELI, B. and MANDELBAUM, A. (2015). Control of Patient Flow in Emergency Departments: Multiclass Queues with Feedback and Deadlines. Forthcoming in Operations Research.
- Hwang, U., McCarthy, M. L., Aronsky, D., Asplin, B., Crane, P. W., Craven, C. K., Epstein, S. K., Fee, C., Handel, D. A., Pines, J. M., Rathlev, N. K., Schafermeyer, R. W., Zwemer, F. L. and Bernstein, S. L. (2011). Measures of Crowding in the Emergency Department: A Systematic Review. *Academic Emergency Medicine* 18 527–538.
- IBRAHIM, R. and WHITT, W. (2011). Wait-Time Predictors for Customer Service Systems with Time-Varying Demand and Capacity. Operations Research 59 1106–1118. MR2864327
- IHI (2011). Patient First: Efficient Patient Flow Management Impact on the ED. Institute for Healthcare Improvement. http://www.ihi.org/knowledge/Pages/ImprovementStories/PatientFirstEfficientPatientFlowManagementED.aspx.
- Janssen, A. J. E. M., van Leeuwaarden, J. S. H. and Zwart, B. (2011). Refining Square-Root Safety Staffing by Expanding Erlang C. *Operations Research* **56** 1512–1522. MR2872017
- JCAHO (2004). JCAHO Requirement: New Leadership Standard on Managing Patient Flow for Hospitals. *Joint Commission Perspectives* 24 13–14.
- Jennings, O. B. and de Véricourt, F. (2008). Dimensioning Large-Scale Membership Services. *Operations Research* **56** 173–187. MR2402225
- Jennings, O. B. and de Véricourt, F. (2011). Nurse Staffing in Medical Units: A Queueing Perspective. *Operations Research* **59** 1320–1331. MR2872002
- JOUINI, O., DALLERY, Y. and AKSIN, O. Z. (2009). Queueing Models for Full-Flexible Multi-class Call Centers with Real-Time Anticipated Delays. *International Journal of Production Economics* 120 389–399.
- Kaplan, R. S. and Porter, M. E. (2011). How to Solve the Cost Crisis in Health Care. Harvard Business Review 89 46–64.
- KC, D. and TERWIESCH, C. (2009). Impact of Workload on Service Time and Patient Safety: An Econometric Analysis of Hospital Operations. *Management Science* 55 1486– 1498.
- Kelly, F. P. (1979). Markov Processes and Reversibility. Wiley.
- Kim, S. H. and Whitt, W. (2014). Are Call Center and Hospital Arrivals Well Modeled by Nonhomogeneous Poisson Processes? *M&SOM* **16** 464–480.
- Koçağa, Y. L., Armony, M. and Ward, A. R. (2015). Staffing Call Centers with Uncertain Arrival Rates and Co-sourcing. *Production and Operations Management* n/a-n/a.
- Leite, S. C. and Fragoso, M. D. (2013). Diffusion Approximation for Signaling Stochastic Networks. *Stochastic Processes and their Applications* **123** 2957–2982. MR3062432
- Long, E. F. and Mathews, K. M. (2012). "Patients Without Patience": A Priority Queuing Simulation Model of the Intensive Care Unit. Working paper.
- MAA, J. (2011). The Waits that Matter. The New England Journal of Medicine 364 2279–2281.
- Maman, S. (2009). Uncertainty in the Demand for Service: The Case of Call Centers and Emergency Departments. Master's thesis, Technion—Israel Institute of Technology.
- Maman, S., Zeltyn, S. and Mandelbaum, A. (2011). Uncertainty in the Demand for Service: The Case of Call Centers and Emergency Departments. Working paper.

- MANDELBAUM, A., MOMCILOVIC, P. and TSEYTLIN, Y. (2012). On Fair Routing from Emergency Departments to Hospital Wards: QED Queues with Heterogeneous Servers. *Management Science* **58** 1273–1291.
- MANDELBAUM, A. and STOLYAR, S. (2004). Scheduling Flexible Servers with Convex Delay Costs: Heavy-Traffic Optimality of the Generalized $c\mu$ -Rule. Operations Research 52 836–855. MR2104141
- MANDELBAUM, A., TROFIMOV, V., GAVAKO, I. and NADJHAHROV, E. (2013). Home-Hospital (Rambam): Readmission Analysis. http://seeserver.iem.technion.ac.il/databases/Docs/HomeHospital_visits_return.pdf.
- MARMOR, Y. N. (2003). Developing a Simulation Tool for Analyzing Emergency Department Performance. Master's thesis, Technion—Israel Institute of Technology.
- MARMOR, Y. N. (2010). Emergency-Departments Simulation in Support of Service-Engineering: Staffing, Design, and Real-Time Tracking. PhD thesis, Technion—Israel Institute of Technology.
- MARMOR, Y. N., GOLANY, B., ISRAELIT, S. and MANDELBAUM, A. (2012). Designing Patient Flow in Emergency Departments. *IIE Transactions on Healthcare Systems Engineering* 2 233–247.
- MARMOR, Y. N., ROHLEDER, T., COOK, D., HUSCHKA, T. and THOMPSON, J. (2013). Recovery Bed Planning in Cardiovascular Surgery: A Simulation Case Study. *Health Care Management Science* **16** 314–327.
- McHugh, M., Van Dyke, K., McClelland, M. and Moss, D. (2011). Improving Patient Flow and Reducing Emergency Department Crowding. *Agency for Healthcare Research and Quality*. http://www.ahrq.gov/research/findings/final-reports/ptflow/index.html.
- NADJHAHROV, E., TROFIMOV, V., GAVAKO, I. and MANDELBAUM, A. (2013). Home-Hospital (Rambam): EDA via SEEStat 3.0 to Reproduce "On Patients Flow in Hospitals". http://ie.technion.ac.il/Labs/Serveng/files/HHD/reproducing_flow %_paper.pdf.
- Nestler, S. (2011). Reproducible (Operations) Research: A Primer on Reproducible Research and Why the O.R. Community Should Care About it. *ORMS Today* 38.
- NGUYEN, V. (1994). The Trouble with Diversity: Fork-Join Networks with Heterogeneous Customer Population. The Annals of Applied Probability 1–25. MR1258171
- PLAMBECK, E., BAYATI, M., ANG, E., KWASNICK, S. and ARATOW, M. (2015). Accurate ED Wait Time Prediction. Working paper, Stanford.
- PLONSKI, O., EFRAT, D., DORBAN, A., DAVID, N., GOLOGORSKY, M., ZAIED, I., MANDELBAUM, A. and RAFAELI, A. (2013). Fairness in Patient Routing: Maternity Ward in Rambam Hospital. Technical report.
- RAMAKRISHNAN, M., SIER, D. and TAYLOR, P. G. (2005). A Two-Time-Scale Model for Hospital Patient Flow. *IMA Journal of Management Mathematics* **16** 197–215. MR2204891
- RAMBAM Rambam Health Care Campus, Haifa, Israel. http://www.rambam.org.il/Home+Page/.
- RAMBAMDATA Rambam Hospital Data Repositories. Technion SEELab, http://seeserver.iem.technion.ac.il/databases/HomeHospital/.
- SAGHAFIAN, S., AUSTIN, G. and TRAUB, S. J. (2014). Operations Research Contributions to Emergency Department Patient Flow Optimization: Review and Research Prospects. Working paper.
- SEELAB SEE Lab, Technion—Israel Institute of Technology. http://ie.technion.ac.il/Labs/Serveng/.

- SEESERVER Server of the Center for Service Enterprise Engineering. http://seeserver.iem.technion.ac.il/see-terminal/.
- SEESTAT SEEStat Documentation, Technion—Israel Institute of Technology. http://ie.technion.ac.il/Labs/Serveng/.
- Senderovich, A., Weidlich, M., Gal, A. and Mandelbaum, A. (2015). Queue Mining for Delay Prediction in Multi-class Service Processes. *Information Systems* n/a–n/a.
- Shi, P., Dai, J. G., Ding, D., Ang, J., Chou, M., Jin, X. and Sim, J. (2013). Patient Flow from Emergency Department to Inpatient Wards: Empirical Observations from a Singaporean Hospital. Working paper.
- SHI, P., CHOU, M. C., DAI, J. G., DING, D. and SIM, J. (2014). Models and Insights for Hospital Inpatient Operations: Time-Dependent ED Boarding Time. *Management Science* **24** 13–14.
- Song, H., Tucker, A. L. and Murrell, K. L. (2015). The Diseconomies of Queue Pooling: An Empirical Investigation of Emergency Department Length of Stay. Forthcoming in Management Science n/a-n/a.
- Stolyar, S. (2005). Optimal Routing in Output-Queued Flexible Server Systems. *Probability in the Engineering and Informational Sciences* **19** 141–189. MR2127332
- Sullivan, S. E. and Baghat, R. S. (1992). Organizational Stress, Job Satisfaction, and Job Performance: Where Do We Go from Here? *Journal of Management* 18 353–375.
- Sun, J. (2006). The Statistical Analysis of Interval-Censored Failure Time Data. Springer. MR2287318
- HCA NORTH TEXAS Hospitals are Moving at the Speed of Life—Real Time Delays Announcement Web-page of Emergency Departments in North Texas, USA. FASTERTX.COM, http://hcanorthtexas.com/.
- Tezcan, T. (2008). Optimal Control of Distributed Parallel Server Systems Under the Halfin and Whitt Regime. *Math of Operations Research* **33** 51–90. MR2393541
- Thompson, S., Nunez, M., Garfinkel, R. and Dean, M. D. (2009). Efficient Short-Term Allocation and Reallocation of Patients to Floors of a Hospital During Demand Surges. *Operations Research* **57** 261–273.
- THORIN, O. (1977). On the Infinite Divisibility of the Lognormal Distribution. Scandinavian Actuarial Journal 1977 121–148. MR0552135
- TSEYTLIN, Y. (2009). Queueing Systems with Heterogeneous Servers: On Fair Routing of Patients in Emergency Departments. Master's thesis, Technion—Israel Institute of Technology.
- TSEYTLIN, Y. and ZVIRAN, A. (2008). Simulation of Patients Routing from an Emergency Department to Internal Wards in Rambam Hospital. OR Graduate project, IE&M, Technion.
- Tukey, J. W. (1977). Exploratory Data Analysis. Addison Wesley.
- MEDICARE USA (2013). Hospital Compare: 30-Day Death and Readmission Measures Data. http://www.medicare.gov/HospitalCompare/Data/RCD/30-day-measures.aspx.
- WARD, A. and Armony, M. (2013). Blind Fair Routing in Large-Scale Service Systems with Heterogeneous Customers and Servers. Operations Research 61 228–243. MR3042753
- WHITT, W. (2012). Fitting Birth-and-Death Queueing Models to Data. Statistics and Probability Letters 82 998–1004. MR2910048
- YOM-TOV, G. B. (2010). Queues in Hospitals: Queueing Networks with ReEntering Customers in the QED Regime. PhD thesis, Technion—Israel Institute of Technology.
- Yom-Tov, G. B. and Mandelbaum, A. (2014). Erlang-R: A Time-Varying Queue with Reentrant Customers, in Support of Healthcare Staffing. *M&SOM* 16 283–299.

ZACHARIAS, C. and Armony, M. (2013). Joint Panel Sizing and Appointment Scheduling in Outpatient Care. Working paper, NYU.

ZAIED, I. (2011). The Offered Load in Fork-Join Networks: Calculations and Applications to Service Engineering of Emergency Department. Master's thesis, Technion—Israel Institute of Technology.

Zeltyn, S., Marmor, Y. N., Mandelbaum, A., Carmeli, B., Greenshpan, O., Mesika, Y., Wasserkrug, S., Vortman, P., Schwartz, D., Moskovitch, K., Tzafrir, S., Basis, F., Shtub, A. and Lauterman, T. (2011). Simulation-Based Models of Emergency Departments: Real-Time Control, Operations Planning and Scenario Analysis. *Transactions on Modeling and Computer Simulation (TOMACS)* 21.

MOR ARMONY
STERN SCHOOL OF BUSINESS, NYU
44 WEST 4TH STREET
NEW YORK, NY 10012
E-MAIL: marmony@stern.nyu.edu

AVISHAI MANDELBAUM
FACULTY OF INDUSTRIAL ENGINEERING
AND MANAGEMENT
TECHNION—ISRAEL INSTITUTE
OF TECHNOLOGY
TECHNION CITY, HAIFA, ISRAEL, 32000
E-MAIL: avim@ie.technion.ac.il

YULIA TSEYTLIN
IBM HAIFA RESEARCH LAB
HAIFA UNIVERSITY CAMPUS
MOUNT CARMEL, HAIFA, ISRAEL, 31905
E-MAIL: yuliatse@gmail.com

SHLOMO ISRAELIT
DIRECTOR, EMERGENCY TRAUMA DEPARTMENT
RAMBAM HEALTH CARE CAMPUS (RHCC)
6 Ha'ALIYA STREET
HAIFA, ISRAEL 31096
E-MAIL: s_israelit@rambam.health.gov.il

Yariv N. Marmor
Department of Industrial Engineering
and Management
ORT Braude College
Karmiel, Israel
Health Care Policy
and Research Department
Mayo Clinic

200 First Street SW Rochester, MN, USA, 55905 E-mail: myariv@braude.ac.il

Galit B. Yom-Tov
Faculty of Industrial Engineering
And Management
Technion—Israel Institute
of Technology
Technion city, Haifa, Israel, 32000
E-mail: gality@tx.technion.ac.il