

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. Some costs are inevitable while others are incurred directly due to inefficient use of resources or indirectly due to suffering patients. All these costs are exacerbated by the stochasticity of hospital systems, which is often manifested by congestion and long delays in patient care. A queueing-network view, of patient flow in hospitals, is thus natural for studying and improving its performance.

The goal of our research is to explore patient flow data through the lenses of a queueing scientist. More specifically, we use exploratory data analysis (EDA) to study patient flow 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 resolutions and operational regimes are relevant for modeling patient length of stay in the Internal Wards (IWs)? Towards fair routing of patients from the ED to the IWs, how is workload measured (via bed occupancy levels or patient turnover rates)? 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 raises 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 their data for universal access, which enables reproducibility of our findings through a user-friendly platform. This will hopefully stimulate readers to carry out their own EDA, to be followed by new models and theory, which would ultimately lead to much needed improvements in hospital patient flow and overall performance.

Keywords and phrases: Queueing Models, Queueing Networks, Healthcare, Patient flow, EDA

CONTENTS

1	Introduction	3
1.1	Patient Flow Focus	4
1.2	EDA, the scientific paradigm and queueing science	5
1.3	Apologies to the Statistician	7
1.4	Rambam hospital	7
1.4.1	The ED+IW network	9
1.4.2	Data Description	9
1.5	Some hints to the literature	9
1.5.1	A proof of concept	11
2	Emergency Department	11
2.1	Basic facts	11
2.2	Exploratory Data Analysis	12
2.2.1	Time dependency	12
2.2.2	Fitting a simple model to a complex reality	12
2.2.3	State dependency	15
2.3	Research Opportunities	17
3	Internal Wards	19
3.1	Basic facts	19
3.2	EDA: LOS—a story of multiple time scales	20
3.2.1	Research Opportunities	22
3.3	EDA: Operational regimes and economies of scale	24
3.3.1	In what regime do IWs operate? Can QED- and ED- regimes co-exist?	25
3.3.2	Research Opportunities	25
3.3.3	Diseconomies of scale (or how ward size affects LOS)	26
3.3.4	Research opportunities	27
4	Transfer from the ED to IWs	27
4.1	Basic facts	28
4.2	Delays in transfer	29
4.2.1	Research Opportunities	30
4.3	Influence of transfer delays on the ED	30
4.3.1	Research Opportunities	31
4.4	Causes of delay	32
4.5	Fairness in the ED-to-IW process	34
4.5.1	Fairness towards patients	34
4.5.2	Research Opportunities	35
4.5.3	Fairness towards staff	35
4.5.4	Research Opportunities	36
5	A system view	37

5.1 **Research opportunities** 38

6 Discussion and concluding remarks 39

6.1 Operational measures as surrogates to overall hospital performance 39

6.2 Multi-dimensional workload 40

6.3 Capacity 41

6.4 Time-scales 42

6.5 Some concluding comments on data-based research—a great opportunity but no less of a challenge 42

6.5.1 Towards a culture of reproducible research in empirical OR/AP 43

Acknowledgements 45

References 47

Appendix: Accessing Data repositories and EDA tools at the SEELab 53

Author’s addresses 54

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; 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 hence present a propitious ground for research in Queueing Theory and, more generally, Applied Probability and Operations Research (OR). Such research would ideally culminate in reduced congestion (crowding) and its accompanying important benefits: clinical, financial, psychological and societal. And a prerequisite for this to happen and for the benefits to accrue, we strongly believe, is that the supporting research is data-based.

Unfortunately, however, operational hospital data is accessible to very few researchers, and patient-level data is in fact publicly unavailable. The reasons span data nonexistence or poor quality, through concerns for patient confidentiality, to proprietorial attitudes of the data owners. We are thus humbly attempting, in this present work, to change this landscape of data-based OR and, in doing so, introduce a new standard. Specifically, we identify and propose research opportunities and challenges that arise from exploratory analysis of ample hospital data. Just as significantly, we also open up our data and make it universally accessible at the Technion IE&M Laboratory for Service Enterprise Engineering ([SEELab](#)): the data

can be either downloaded or analyzed online, through a user friendly platform ([SEESat](#)) for Exploratory Data Analysis (EDA). Our goal is thus to provide an entry to and accelerate the learning of data-based OR of hospitals; Interested researchers can reproduce our EDA, and use it 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; and 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 measure, and that they inherently serve as “surrogates” for other quality of care measures. For example, the rate of readmission to the hospital, within a relatively short time, often serves as a proxy of clinical quality of care. Similarly, the rate of LWBS (Left without being seen) is a common proxy for accessibility to care. In parallel, patient flow has caught the attention of researchers in Operations Research, Applied Probability, Service Engineering and Operations Management, with Queueing Theory serving as a common central thread that connects these disciplines. This is not surprising: hospital systems, being congestion-prone, naturally fit the framework of Queueing Theory, which captures the tradeoffs between (operational) service quality vs. resource efficiency.

Our starting point is that a queueing network encapsulates the operational dimensions of patient flow in hospitals, with the medical units being the nodes of the network, patients are the customers, while beds, medical staff and medical equipment are the servers. But what are the special features of this queueing network in terms of its system primitives, key performance measures and available controls? To address this question, we study an extensive data set of patient flow through the lense 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, as highlighted by our title (“On 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. Furthermore, having

calculated a 90×90 transition matrix between hospital wards (see Figure 2 in [EV](#)), we shall focus on the sub-system we call the ED+IW network (more on that momentarily). To be concrete, we analyze patient flow data of the Emergency Department (ED) ([§2](#)), Internal Wards (IWs) ([§3](#)), and the transfer of patients from the ED to IWs ([§4](#)). [Section 5](#) discusses the interplay between these three components, and thus provides an integrative view of the system as a whole. We offer final commentary in [§6](#), where we also provide a broader discussion of some common themes that arise throughout the paper. Finally, we provide data access instruction and documentation, as well as EDA logistics in the Appendix. We encourage interested readers to refer to [EV](#), which is an extended version of the present paper: it provides more elaborate discussions of various issues and it covers additional topics, not included here due to space and focus considerations.

1.2. *EDA, the scientific paradigm and queueing science.* Our approach of learning from data is in the spirit of Tukey’s Exploratory Data Analysis (EDA) ([Tukey, 1977](#)). To quote from [Brillinger \(2002\)](#), John Wilder Tukey “...recognized two types of data analysis: exploratory data analysis (EDA) and confirmatory data analysis (CDA). In the former the data are sacred while in the latter the model is sacred. In EDA the principal aim is to see what the data is “saying”. It is used to look for unexpected patterns in data. In CDA one is trying to disconfirm a previously identified indication, hopefully doing this on fresh data. It is used to decide whether data confirm hypotheses the study was designed to test”.

Within the framework of the EDA-CDA dichotomy, here we focus on EDA. This prepares the ground for future CDA which, in the present context, would be the application of data-based (queueing and statistical) models to confirm or refute prevalent hypotheses. Confirmation would contribute insight that supports the management of the originating system(s) (patient flow in hospitals, in this paper); refutation gives rise to new hypotheses, further EDA then CDA, with ultimately new insightful models. This EDA-CDA cycle is routine in natural sciences, where it is commonly referred to as The Scientific Paradigm ([The OCR Project \(IBM-Rambam-Technion\), 2011](#)). Human complexity forced the paradigm onto Transportation Science ([Herman, 1992](#)) and Behavioral Economics ([Camerer, Loewenstein and Rabin, 2003](#)), and the present study aims at doing the same for the analysis of patient flow in hospitals. A similar approach has already proved successful in other operational settings, including semi-conductor manufacturing ([Chen et al., 1988](#)), telecommunication ([Leland et al., 1994](#)), new product development ([Adler et al., 1995](#)) and, most recently, call centers (see [Man-](#)

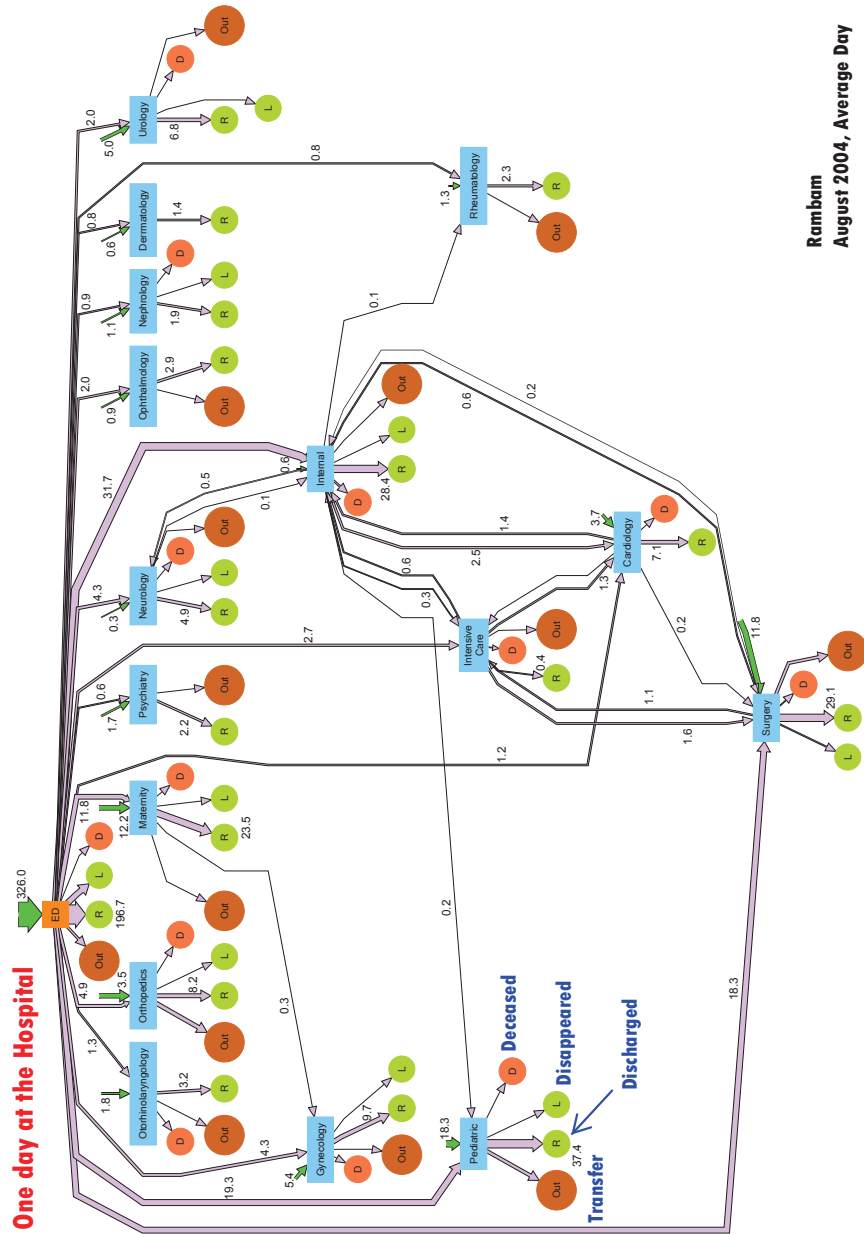


FIG 1. Patient Flow (Process Map) at inter-ward resolution. For example, during the period over which the flow was calculated (August 2004), 326 patients arrived to the ED per day on average, and 18.3 transferred from the ED to Surgery. (To avoid clutter, arcs with monthly flow below 4 patients were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

delbaum, Sakov and Zeltyn (2000) and Brown et al. (2005) for the empirical findings, and Gans, Koole and Mandelbaum (2003) and Aksin, Armony and Mehrotra (2007) for surveys on follow-up work).

1.3. *Apologies to the Statistician.* We conclude our EDA discussion with two “apologies” to the Statistician. Firstly, the goals of the present study, its target audience and space considerations render secondary the role of “rigorous” statistical analysis (e.g. hypothesis testing, confidence intervals, model selection). Indeed, we believe that its intentional omission is both necessary and justified. Accordingly, we either mention a statistical rationalization only in passing (e.g. for the fact that patients Length-of-Stay (LOS) distribution is Log-Normal in resolution of days and it is a mixture of Normal distributions in an hourly resolution), or we simply content ourselves with a convincing visual evidence (the privilege of having a large data set).

Secondly, our data originates from a single Israeli hospital, operating during 2004–2008. This raises doubts regarding the generality of the scientific and practical relevance of the present findings, and rightly so. Nevertheless, other studies of Israeli hospitals (Marmor (2003); Tseytlin (2009) and EV) indicate that these hospitals have many common features. This still leaves the concern that Israeli hospitals during the considered period are perhaps too “unique”—however, in many significant ways they are not: hospitals are slow to change and, more concretely, a parallel recent study by Shi et al. (2012) in a major Singapore hospital, together with other privately-communicated empirical research by colleagues, reveal phenomena that are common across hospitals (e.g. the LOS distributions in Figure 9). All in all, our hope is that reading the manuscript will dispel all doubts concerning its broad relevance and significance (practical, statistical and scientific in general).

1.4. *Rambam hospital.* Our data originates at the 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), and at the resolution level of Figure 1. In particular, 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.

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 lead to a lack of focus. Instead, and although our data encompasses the entire hospital, we chose to focus on a sub-network that consists of the main ED (adult Internal, Orthopedics, Surgery, and Trauma) and five IWs, denoted by A through E; see Figure 2. This sub-network, referred to

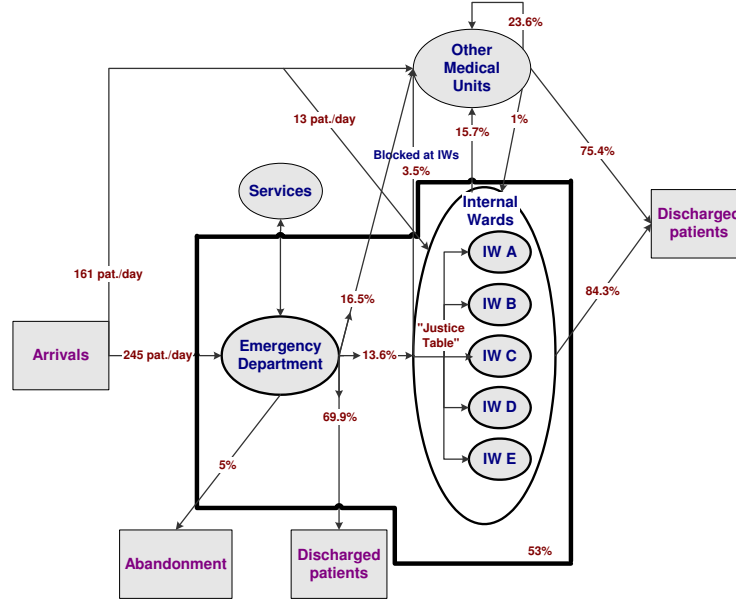


FIG 2. The ED+IW sub-network

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 Figure 2 in EV, approximately 53% of the patients entering the hospital remain within this sub-network, and 21% of those are hospitalized in the IWs; indeed, the network is fairly isolated in the sense that its interactions with the rest of the hospital are minimal. To wit, virtually all arrivals into the ED are from outside the hospital, and 93.5% of the patient transfers into the IWs are either from outside the hospital or from within the ED+IW network.

While focusing on the ED+IW network, we nevertheless reap the benefits of having access to overall hospital data. One such benefit is the use of other hospital units as reference points—this enhances understanding of specific phenomena that arise from the ED+IW data.

1.4.1. *The ED+IW network.* The main ED has 40 beds and it treats on average 245 patients daily. An internal patient, whom an ED physician decides to hospitalize, is directed to one of the five Internal wards. The IWs have about 170 beds that accommodate around 1000 patients per month. Internal Wards are responsible for the treatment of a wide range of internal conditions, thus providing inpatient medical care to thousands of patients each year. Wards A–D share more or less the same medical capabilities—each can treat similar (multiple) types of patients. Ward E, on the other hand, attends to only the less severe (walking) cases; in particular, this ward cannot admit ventilated patients.

1.4.2. *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, some demographic data, and more. 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, which has enabled the presently reported EDA and more.

1.5. *Some hints to the literature.* Patient flow in hospitals has been studied extensively. Readers are referred to the many papers in [Hall \(2006\)](#) and the recent [Shi et al. \(2012\)](#)—both providing leads to further references. In the present subsection, we merely touch on published work, along the three dimensions that are most relevant for 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 further cited as we go along. This subsection concludes with what can be viewed as “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 ED flows in within. 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. operating-room schedules that are customized to physician needs rather than being operationally optimized). 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 in-patients; and 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. \(2012\)](#) 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. Two recent handbooks on System Scheduling and OM in Healthcare are [Hall \(2012\)](#) and [Denton \(2013\)](#)—both include chapters worth reading and additional leads on OR/OM and queueing perspectives of patient flow. Of special interest is Chapter 8 in [Hall \(2012\)](#), 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 \(2007\)](#) and [Yom-Tov and Mandelbaum \(2011\)](#) develop (time-varying) queueing networks to 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 \(2009\)](#) use econometric methods to investigate the influence of workload on service time and readmission probability in Intensive Care Units (ICUs).

This inspired [Chan, Yom-Tov and Escobar \(2011\)](#) to model an ICU as a state-dependent queueing network, in order to gain insight on how speedup and readmission effects influence the ICU.

1.5.1. *A proof of concept.* The present research has already provided the empirical foundations for several graduate theses, each culminating in one or several 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 customer returns to the ED ([Yom-Tov and Mandelbaum, 2011](#)) 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 \(2013\)](#) develops scheduling controls that help ED physicians choose between newly-arriving vs. in-process patients, while still adhering to triage constraints ([Huang, Carmeli and Mandelbaum, 2011](#)). These are all examples of the EDA-CDA process, alluded to in Subsection 1.2.

2. Emergency Department. Patient flow in the Emergency Department (ED) is a complex process that involves a multitude of interrelated steps (see Figure 8 of [EV](#)). This process has been widely investigated, both academically ([Hall, 2006](#)) and in practice ([IHI, 2011](#); [McHugh et al., 2011](#)). We shall hence be content here with its empirical macro view, which already turns out to be highly informative. Specifically, we view the ED as a black-box, and then highlight interesting phenomena that relate to its patient arrivals, departures, and occupancy counts. Our EDA underscores the importance of including time- and state- dependent effects in a queueing model of the ED. Yet, and albeit this dependence, it also reveals that a simple stationary model may well fit patient-count during periods when the ED is most congested. For limited purposes, therefore, our EDA supports the use of “black-box” stationary models for the ED, which has been prevalent in the literature (e.g. [Green et al. \(2006\)](#) and [de Bruin et al. \(2009\)](#)).

2.1. *Basic facts.* Rambam’s main ED attends to 200–250 patients daily: close to 60% are classified as Internal (general) patients and the rest are Surgical/Orthopedic, excluding a few per day that suffer from multiple trauma. The ED has three major areas: Trauma acute, Internal acute, and Surgical/Orthopedic acute; some of the patients in the latter two are “Walking” patients that do not require a bed. While there are formally 40 beds in the ED, this bed capacity is highly flexible and can be doubled and more. Hence there is effectively no upper bound on how many patients can simultaneously

reside within the ED—either in beds or sitting and waiting. 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. Moreover, within the ED, we focus on Internal (general) patients, in beds or walking: they both constitute the majority of ED patients and give rise to most 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 4.2). 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 throughout the process and are not LWBS or LAMA). Finally, 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.

2.2. *Exploratory Data Analysis.* In this section we highlight some of our EDA findings that relate to patient arrivals and patient-count distribution. We observe both time- and state-dependent behavior of these entities, some of which are not readily explained by existing queueing models.

2.2.1. *Time dependency.* As observed also in [Green, Kolesar and Whitt \(2007\)](#), the ED hourly arrival rate varies significantly during the day. In Rambam’s ED, it varies by a factor of almost 10; See Figure 3. We also observe a time-lag between the arrival rate and occupancy levels, which is due to the former changing significantly during a patient LOS ([Bertsimas and Mourtzinou, 1997](#)). This lag must be accounted for in staffing recommendations ([Feldman et al., 2008](#); [Green, Kolesar and Whitt, 2007](#)).

Analyzing the same data, [Maman \(2009\)](#) also found support for the daily arrival process to fit a time-varying Poisson process, but with heterogeneity levels across days such that the *arrival rate* itself must be *random* (slightly over-dispersed). [Kim and Whitt \(2013\)](#) identified similar patterns in a large Korean hospital. The time-varying arrivals contribute to an overall time varying ED environment, which we focus on next.

2.2.2. *Fitting a simple model to a complex reality.* Figure 4 (left) shows 24 patient-count histograms for internal ED patients, each corresponding to

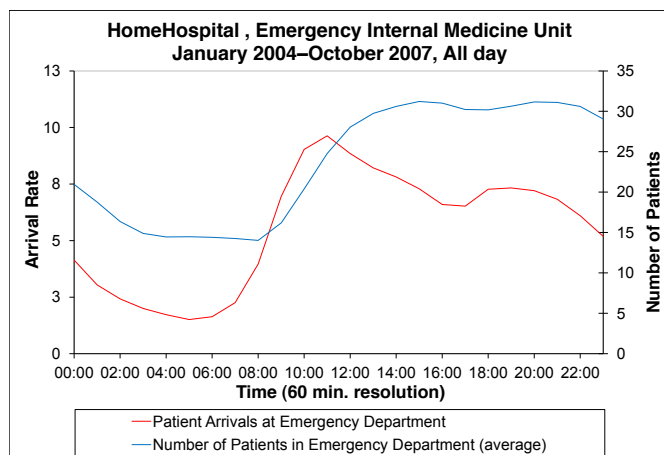


FIG 3. Average number of patients and arrival rate by hour of the day

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) during the AM (3–9AM), and high (30 patients) during the PM (12–11PM); with two transitional periods of low-to-high (9AM–12PM) and high-to-low (11PM–3AM). We refer to these four periods as the four “occupancy regimes”.

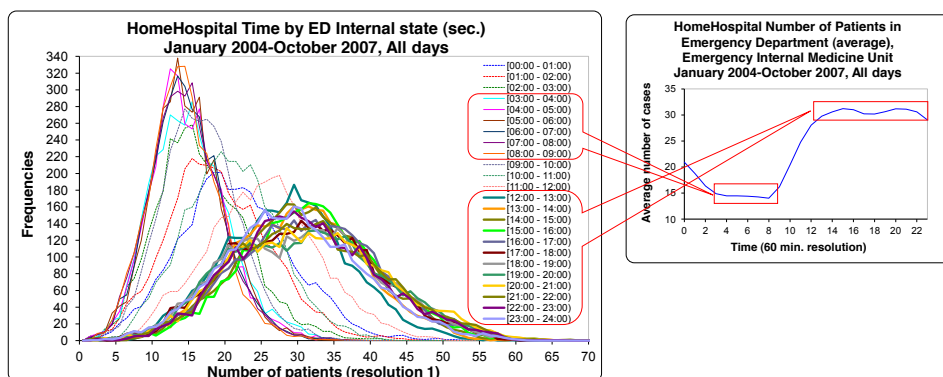


FIG 4. Internal ED Occupancy histogram by hour of the day

Interestingly, when asking SEESat to fit a mixture of three normal dis-

tributions to the ED occupancy distribution, the fit algorithm *automatically* detects the low, high and transitional phases (See Figure 5).

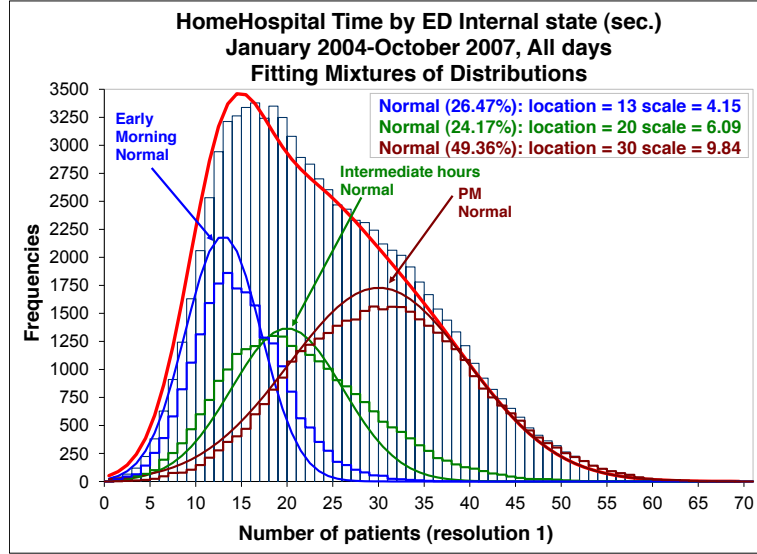


FIG 5. Fitting a mixture of three normal distributions to the ED occupancy distribution

Further EDA (in EV) suggests 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, it may also be described as an $M/M/N+M$ model with equal service and abandonment (LWBS, LAMA, or Absconded) rates. It follows that one cannot conclusively select a model through its empirical steady-state distribution—which is a trap that is easy to fall into and from which Whitt (2012) saved us.

One is thus led to the relevance-limits of “black-box” ED models: they may support operational decisions that depend only on total patient count but not on (and neither do 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”, then its reversibility accommodates also applications that *do change* total count: for example, ambulance diversion in the face of total count that exceeds a certain threshold, which then truncates the count to this threshold (and the steady-state distribution is truncated corre-

spondingly; see Kelly (1979)). On the other hand, black-box models cannot support ED staffing (e.g. Yom-Tov and Mandelbaum (2011) acknowledges some internal network dynamics), or ambulance diversion that depends on the number of boarding patients (awaiting hospitalization). We discuss this further in Section 2.3.

2.2.3. *State dependency.* In addition to time-dependent effects, we discover 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 (in bed or walking): the left graph displays the total service rate, and the right graph shows the rate per Internal patient. These graphs cannot arise from commonly used (birth-death) queueing models such as $M/M/N$ (total rate that is linearly increasing up to a certain point and then constant) or $M/M/\infty$ (constant per patient). In contrast, the per-patient service rate has an interval ($11 \leq L \leq 20$) where it is increasing in L , which is between two intervals of decrease. (The noise at the extremes, $L \leq 3$ and $L \geq 55$, is due to small sample sizes.) (Note that Batt and Terwiesch (2012) and Kc and Terwiesch (2009) also found evidence for a load dependent service rate.)

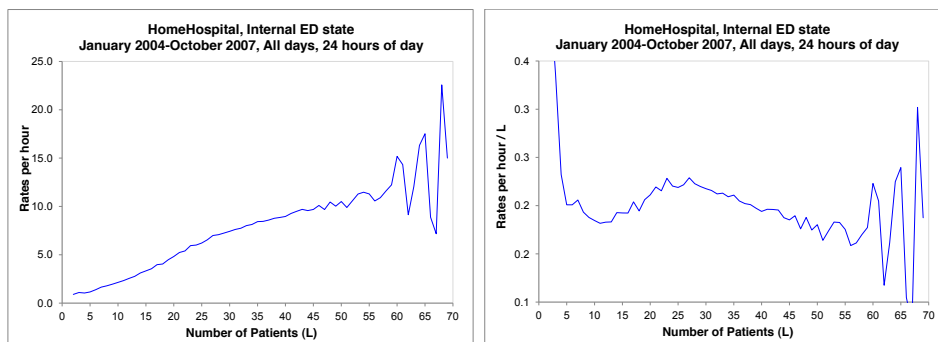


FIG 6. *Service rate and service rate per patient as a function of L*

What can cause this particular state-dependence of the service rate per patient? We start with the “slowdown” ($L \geq 25$) which, in a congested ED, is to be expected under any of the following scenarios:

- *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 slow down (Sullivan and Baghat, 1992).
- *Choking*: Service slowdown may also be attributed to so-called resource “choking”: medical staff becomes increasingly occupied with caring for to-be-transferred (boarding) ED patients (who create work while they wait and, moreover, their condition could actually deteriorate), that they might end up choking off the throughput of the to-be-released patients (see Figure 13 in Section 4.3). 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, slowdown as well as speedup may be attributed to the combination of time dependent arrivals and heterogenous patient mix (Marmor et al., 2011). We now expand on the speedup effect.

As opposed to the slowdown, the apparent speedup ($10 \leq L \leq 25$) turns out to be an artifact of biased sampling due to patient-heterogeneity and time-variability (as observed in Section 2.2.1). 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.

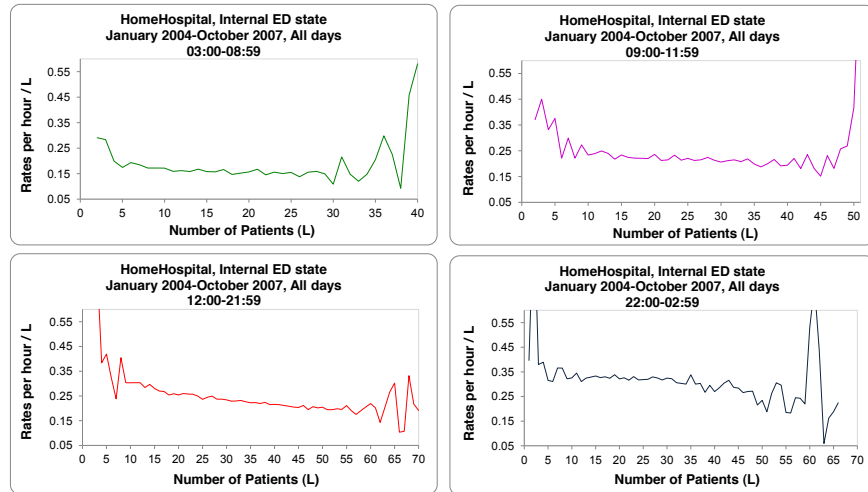


FIG 7. Service rate per patient as a function of L by occupancy regime

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. So what appears to be a speedup (increasing graph) is merely a weighted average of non-increasing graphs with state-dependent weights.

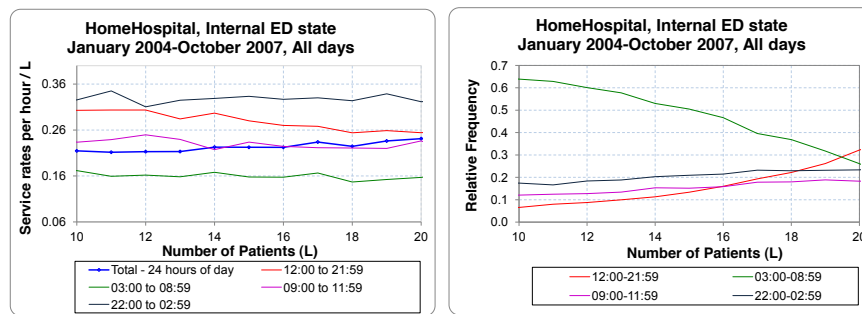


FIG 8. Service rate as a function of $10 \leq L \leq 20$ and Relative frequency (weight) of occupancy regime per L

2.3. Research Opportunities. Our EDA leaves many open questions for further exploration. For example: What is the source for the unique shape of the time dependent arrival pattern, which is common in many service systems (including hospitals across the globe and call centers)? What is the dominant factor in explaining the slowdown in service rate per patient, and what can be done to alleviate this slowdown? Considering time and state dependency, how does one separate these two effects? Which one is more dominant and under what circumstances?

The observations in this section raise some broader research directions within a few (somewhat overlapping) model dimensions: granularity, performance metrics, and applications.

- *Model granularity:* Our focus in this section has been on overall ED (Internal) patient count. This aggregates ED dynamics into merely arrivals and departures which, as described in Subsection 2.2.2, yields a useful black-box model but of a limited scope. In contrast, one could consider a micro-level that acknowledges explicitly all events within the ED, for example events arising during patient routes, service times and resources utilization—more on that momentarily.

The macro- and micro-models are two extreme cases of model granularity, with a range of levels in between (e.g. Yom-Tov and Mandelbaum (2011), Huang, Carmeli and Mandelbaum (2011), and Mar-

mor et al. (2012)). In addition, one may also use a combination of the macro-micro models where various black boxes serve as analytical nodes within an elaborate system.

The granularity level to be used depends on target application, availability of data and analytical techniques. Put differently, what kind of operational decisions require which model granularity? And more concretely, for our proposed “blackbox” model—what are the purposes for which this simple model is sufficient and what are this model’s limitations?

- *Performance metrics:* There are numerous ED performance metrics that have not been discussed here or have merely been touched upon. These pertain to time-till-first-consultation, length of stay (LOS), causes of operational delays (scarce resources, synchronization gaps), abandonment types (LWBS, LAMA, absconded), readmissions, workloads and offered load, bed utilization, boarding times, staff-to-bed ratios, and customers who are blocked upon their ED arrival (e.g. ambulance diversion) or departure (boarding). Of special importance are ED congestion metrics (Hwang et al. (2011) lists over 70), which have given rise to prevalent crowding indices (e.g. Bernstein et al. (2003); Hoot et al. (2007)).

Metrics must be assigned units, be measured (often a challenge) or be estimated. As an example for the latter, (im)patience of LWBS patients (i.e. the time these patients are willing to wait before leaving) is observed only if patients announce their departure. Otherwise patients are either served, in which case their waiting time provides a lower bound for their (im)patience, or they are discovered missing when called for service, which provides an upper bound. Statistical inference of ED (im)patience therefore requires novel models and methods: these would combine current-status (Sun, 2006) and survival-analysis (Brown et al., 2005) setups—in the latter, abandonment times are observed, while they are not in the former.

- *Applications:* Applications of queueing models for ED patient flow include the following categories: ED design, capacity sizing, staffing (e.g., Yom-Tov and Mandelbaum (2011)), and flow control (e.g., Allon, Deo and Lin (2010); Dobson, Tezcan and Tilson (2013); Hagtvedt et al. (2009); Huang, Carmeli and Mandelbaum (2011)).

An outcome of ED design is flow architecture (Marmor et al., 2012). Related examples that would enjoy research are operational (fast-track) vs. clinical priorities (see also Zeltyn et al. (2011)), physician-led triage vs. the prevalent nurse-led (Burström et al., 2012; Oredsson

et al., 2011), and the creation of a dedicated sub-ED (e.g. for patients with chest-pain; Zalenski et al. (1998)).

An important broader challenge is the evaluation of ED (Emergency Department) vs. ER (Emergency Room) designs: the former functions as a bonafide hospital ward that provides treatment to most patients, while the latter is mainly a router to hospital wards, which treats scarcely few. A final challenge is to model information within the ED. Indeed, ED processes are geared towards accumulation of information, up to a level that suffices to support the final decision of discharge vs. admit; and the tradeoff in these processes is the classical exploration (e.g. administer new tests or additional treatment) vs. exploitation (i.e. make a final decision based on the current information); see Gittins, Glazebrook and Weber (2011) and Huang, Carmeli and Mandelbaum (2011).

3. Internal Wards. Internal Wards (IWs), often referred to as General Internal Wards or Internal Medicine Wards, are the “clinical heart” of a 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 resources (beds, physicians). These characteristics are attributed to IW inflow design, capacity management and operational policies (e.g. discharge procedures, physician rounds).

3.1. *Basic facts.* Rambam hospital has five Internal wards. Wards A–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 IWs. For example, bed capacity ranges from 24 to 45 beds and Average LOS (ALOS) from 3.9 to 6 days.

IWs B and E are by far the smallest (least number of beds) and the “fastest” (shortest ALOS, highest throughput). The superior operational performance of 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 qual-

TABLE 1
Internal wards operational profile

	Ward A	Ward B	Ward C	Ward D	Ward E
Average LOS (days) (STD)	6.0 (7.9)	3.9 (5.4)	4.9 (10.1)	5.1 (6.6)	3.7 (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) capacity (# beds)	45 (52)	30 (35)	44 (46)	42 (44)	24
Mean # patients per bed per month	4.58	6.45	4.77	5.16	7.44
Readmission rate (within 1 month)	10.6%	11.2%	11.8%	9.0%	6.4%

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).

ity of care. To get a grasp on that quality, we use its operational (accessible hence common) proxy, namely patient readmission rate (proportion of 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 wards. 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 later in Section 3.3.3.

3.2. *EDA: LOS—a story of multiple time scales.* In this section, we examine the distribution of the 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. It does so also in other service systems though an explanation for its prevalence is still lacking ([Brown et al., 2005](#)). When considering an hourly resolution, however, a completely different distribution shape is observed: there are peaks that are periodically 24 hours apart, which corresponds to a *mixture* of daily distributions. (We found that a normal mixture fits usefully well, as depicted by the 7 normal mixture-components over the range of 0–150 hours in Figure 9.)

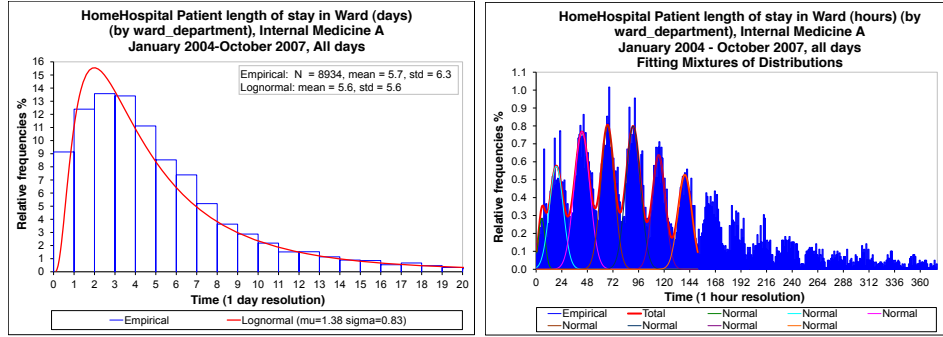


FIG 9. LOS distribution of IW A in two time-scales: daily and hourly

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 discharge letters (after finishing the morning rounds); then nurses take care of paperwork, instructing patients (and their families) 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 yields an explanation for the observed peaks that are spaced 24 hours apart. The variation around these peaks is determined by the arrival process: patients are hospitalized in IWs almost exclusively over a 12-hour period (10am–10pm), with a peak in arrival rate between 3pm–7pm (Figure 10). Following similar observations in a Singaporean hospital, Shi et al. (2012) offer a 2-time-scale mathematical model that supports our EDA.

Note that the arrival process to the IWs is mostly a departure process from the ED, and hence the timing of its peak (3pm–7pm) is naturally coupled with IW discharge peaks (3pm–4pm). In other words, and as further discussed in Section 5, the discharge policy from IWs significantly influences ED congestion. This led Shi et al. (2012) to propose flow-stabilization as a means for reducing ED congestion, which is caused by a ward discharge

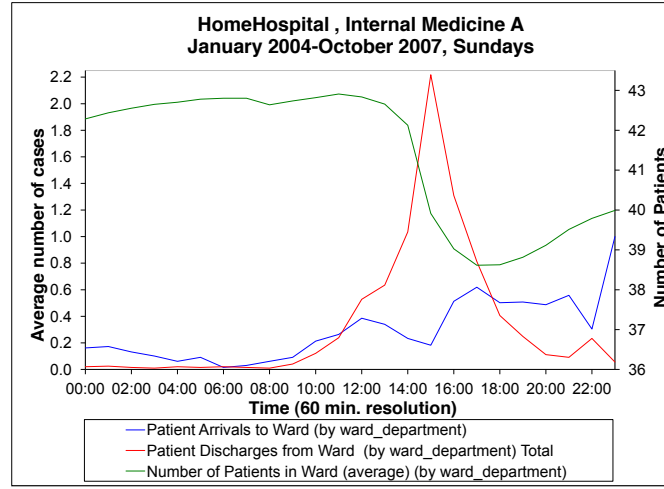


FIG 10. Arrivals, departures, and average number of patients in Internal wards by hour of day

policy that resembles ours.

3.2.1. Research Opportunities. We discuss here multiple time-scales, server identification, workload characterization, protocol mining via LOS distributions, flow control and why Log-Normal.

Multiple Time Scales: 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. 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. Here service times are measured in minutes or parts of an hour, and offered load (workload) is calculated (from arrival and service data) in units of, say, hours of work that arrives per hour of the day.

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 offered load results from both routine and special care, and it varies during the day for at least two reasons (see Equation (1) in [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 6.2.

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., 2012). Or differences in shape of LOS distribution between two Maternity wards (§4.2.1 in EV), which result from differing patient mix, 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, the latter being the drivers of operational performance. This leads to open data-based questions in two directions: either incorporating protocols (e.g. patient priorities, resource scheduling) in queueing models and validating the theoretical LOS distribution against data (performance); or, 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. (2012). We return to it in Section 5.

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? It is 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 in-

finitely divisible, being an integral against a Gamma process ([Thorin, 1977](#)). Can these properties help one explain the empirical Log-Normal service time distribution?

3.3. *EDA: Operational regimes and economies of scale.* An asymptotic theory of many-server queues has been developed in recent years ([Gans, Koole and Mandelbaum \(2003\)](#) can serve as a starting point), 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. For example, in well-run call centers, server utilization could exceed 90% while, at the same time, about half of the customers are served without delay, and those delayed wait one order of magnitude less than their service duration (seconds vs. minutes). The QED regime also exhibits economies of scale in the sense that, as the system grows, operational performance improves (e.g. less wait and less abandonment, under equal workload per server).

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 Leeuwen and Zwart, 2011](#)) and, significantly, it renders QED theory relevant to healthcare systems ([Jennings and de Véricourt, 2011](#); [Yom-Tov and Mandelbaum, 2011](#)). 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 is 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 (§3.3.3) the existence of the economies-of-scale phenomenon in the hospital environment. We shall argue that, although IW beds plausibly operate in the QED regime, there is

nevertheless evidence for diseconomies of scale.

3.3.1. *In what regime do IWs operate? Can QED- and ED-regimes co-exist?* We start by identifying the operational regimes that are relevant to our system of IWs. This system has multiple types of servers (beds, nurses, physicians, medical equipment), and each must be considered separately. 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 one 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 on 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 on call 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 4.2). Performance of physicians is therefore Efficiency Driven.

3.3.2. **Research Opportunities.** We identified two operational regimes, QED and ED, that coexist within the ED+IW: waiting in the ED for IW service. What queuing 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 ED resource (physicians)? Thus, can one conclude that

adding physicians is necessary for reducing transfer delays, while adding beds would have only a marginal impact on these delays? (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.) How would a change of 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 physicians operate in the ED-regime eliminate the economies of scale that one expects to find in QED systems? Empirical observations that will now be presented suggest that this might indeed be the case.

3.3.3. *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 attributed 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 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.

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 number 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) vs. operational efficiency (flexible). This tradeoff has been addressed in queueing models (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

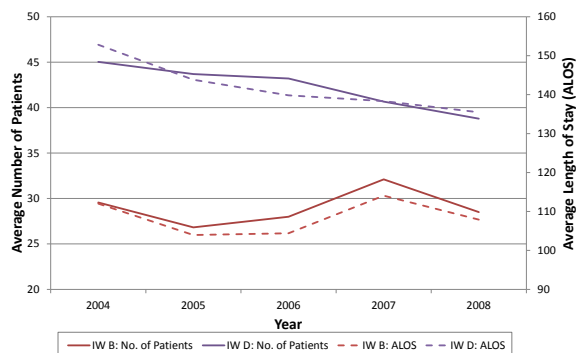


FIG 11. Average LOS and number of patients in Internal wards B and D by year

will surely add novel idiosyncracies.

- *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 (physicians' rounds) and discuss courses of treatment. This is essential as Rambam hospital 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)—this could lead to a prolonged ALOS.

3.3.4. 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 Theory is well equipped to address the operational dimensions of such decisions, it will have to collaborate with other disciplines such as organizational behavior for complete comprehension. 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 our next subject.

4. Transfer from the ED to IWs. The “ED-to-IW” process is the channel of hospitalization for 21% of the Internal ED patients. We focus on

a problem that commonly arises in this process—long patient waiting times in the ED for a transfer to the IWs—and discuss its influence on patients and staff. We view the process in the context of *flow control*, where patients are *routed* from the Emergency Department to Internal Wards.

Routing in *hospitals* differs from other service systems for various reasons including incentive schemes, customers’ (patients’) limited control (often bordering on 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. Studying the transfer process in this setting leads to many research opportunities.

4.1. *Basic facts.* We begin with a short description of the patient transfer process from the ED to the IWs in Rambam hospital. A patient, whom an ED physician decides to hospitalize in an IW, is assigned to one of the 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 Head nurse of the hospital. Once the assigned ward is set, the ward staff prepares for this patient’s arrival. In order for the transfer to start, a bed and medical staff must be available, and the bed and equipment must be prepared for the patient (including potential rearrangement of current IW patients). Up to that point, the 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 responsibility while medical treatment is still obtained from an IW physician.

An integral component of the transfer process is a *routing policy*, or patients assignment algorithm. As described in Section 3.2, Wards A–D provide similar medical services, while Ward E treats only the less severe “walking” 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 patient load among wards, so that staff workload on will be balanced. A survey among 5 additional hospitals in Israel (EV, Section 5.6) reveals that a cyclical routing

policy is very common; yet, some hospitals apply alternative assignment schemes. For example, one hospital uses random assignment by patient ID. Surprisingly, only one of the surveyed hospitals uses an assignment that takes into account real-time bed occupancy.

4.2. *Delays in transfer.* As is customary elsewhere, the operational goal of our hospital is to admit ED 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, 23% of the patients were delayed by more than 4 hours.

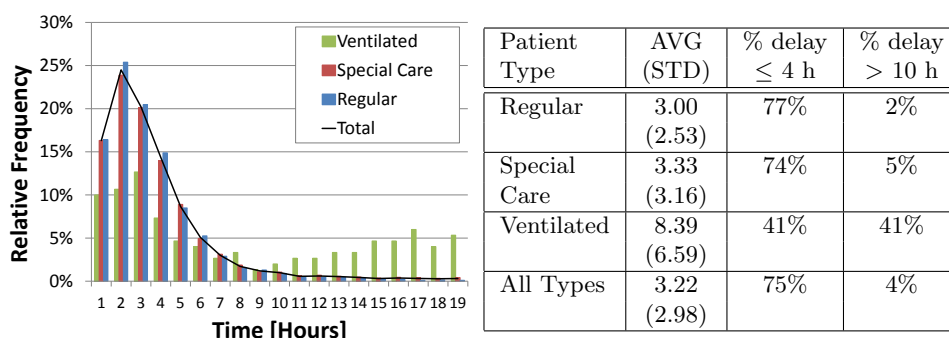


FIG 12. *Transfer time by patient type, in hours*

* Data refer to period 5/1/06–10/30/08 (excluding the months 1–3/07 when Ward B was in charge of an additional sub-ward)

An interesting phenomenon is observed when analyzing transfer delays by patient types. 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 by more than 10 hours. Ventilated patients must have the highest priority in transfer but, in reality, many do not benefit from it.

How can it be that many of the ventilated patients experience such long delays? We observe that the ventilated shorter delays (< 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) (Zhu, Armony and Chan, 2013). Each such TCU has only 4–5 beds. The average occupancy rate of the TCUs at Rambam hospital is 98.6%; the combination of high occupancy

with a small number of beds results in much longer waits during overloaded periods. (b) Ventilated patients require a highly qualified staff to transfer them to their ward (especially since they are attached to an oxygen source). Coordinating such transfers takes longer.

4.2.1. Research Opportunities. Delays in transfer provide additional opportunities to those discussed at the end of §3.2.1. First there is the challenge of deciphering protocols—here ED-to-IW routing—from data such as in Figure 12. Then one would like to be able to analyze and optimize patient-flow protocols in queueing models, specifically here fork-join networks 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).

4.3. Influence of transfer delays on the ED. Patients awaiting transfer (boarding patients) do overload the ED: beds remain occupied while new patients continue to arrive, and the ED staff remains responsible for those 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 may suffer from transfer delays.

Boarding patients may experience significant discomfort while waiting: the ED is noisy, it is not private and does not serve hot meals. In addition, ED patients do not enjoy the best professional medical *treatment* for their particular situation, 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 transfer patients can be significant. Figure 13 shows our estimates of the fraction of time that ED physicians spent caring for the transfer patients, assuming (the Rambam experience) that every such patient requires 1.5 minutes of physician’s time every 15 minutes. We observe that transfer 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*: a phenomenon that is well acknowledged in transportation (Chen, Jia and

Varaiya, 2001) and telecommunication (Gerla and Kleinrock, 1980) and discussed earlier in Section 2.2.3.

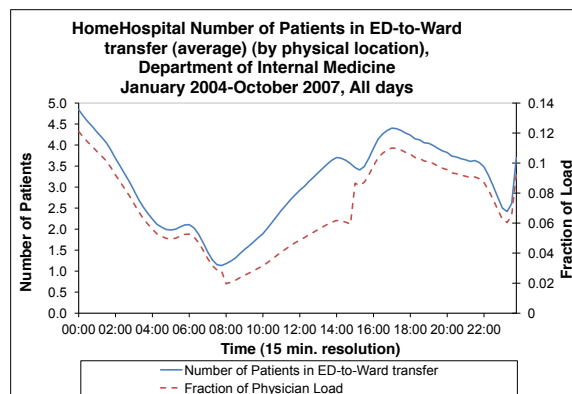


FIG 13. Number of patients in ED-to-IW transfer (A–E) and the fraction of time that ED physicians devote to these patients

To summarize, by improving patient flow from the ED to the IWs, in particular reducing transfer times, hospitals can improve the service and treatment provided to both transfer and in-process patients. In turn, reducing the workload on the ED will lead to a better response to arriving patients and is likely to save lives.

4.3.1. Research Opportunities. The delays in transfer give rise to interesting research questions. For example:

1. *Modeling transfer queue:* Transfer patients may be viewed as customers waiting in queue to be served in the IW. Traditionally, it has been assumed that the customers receive service only once they reach the service station, 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, transfer 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.
2. *Emergency Department architecture:* As described, ED staff takes care of two types of patients: transfer and in-process patients. Each type has its own service requirements, leading to differing service distributions and differing distribution of time between successive treatments.

While transfer patients receive periodic service according to a nearly-deterministic schedule (unless complications arise), in-process service is more random.

One may consider two options for ED architecture: (a) treat transfer and in-process patients together in the same physical location, as is done at Rambam, or (b) move the transfer 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 naturally arises: Under what conditions in each of these ED architectures more appropriate?

Note that the Singapore hospital architecture is even more complicated than (b), as the responsibility for the transfer patients is handed over to IW physicians after a two-hour delay. 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 discuss how different architectures are related to incentive schemes and, in turn, influence delay times.

4.4. *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 and motion study, and further explored our data. We have learned that delays are not only caused by bed unavailability; patients often wait even when there are available beds. Indeed, our data show 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. 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 causes for delay, which are summarized in the Cause-and-Effect (Fishbone) diagram depicted in [Figure 14](#). We elaborate here on two that have interesting modeling aspects.

1. *Input-queued vs. Output-queued system:* Recall that the preparation for a *particular* transfer 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 in which patients are placed in a “common” queue, are routed

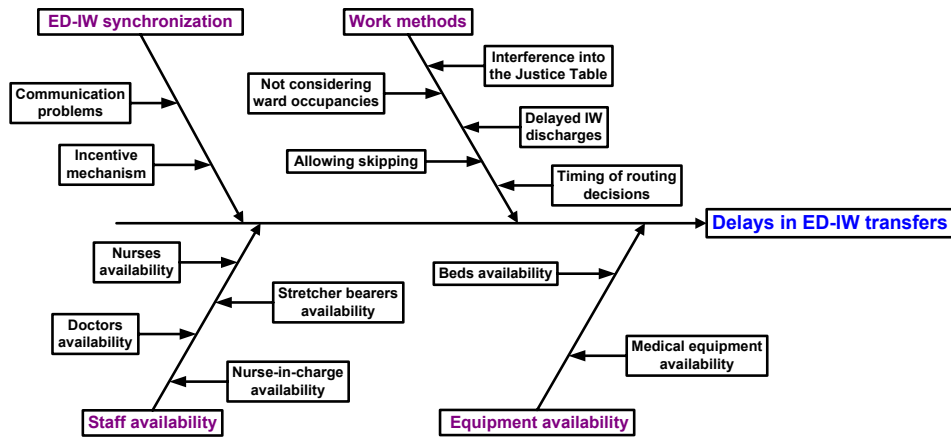


FIG 14. ED-to-IW delays: Causes and effects chart

to an IW only once at the head of the queue and *any* of the beds in the IWs become available. The latter is referred to as an *input-queued* scheme. Figure 15 depicts the two schemes.

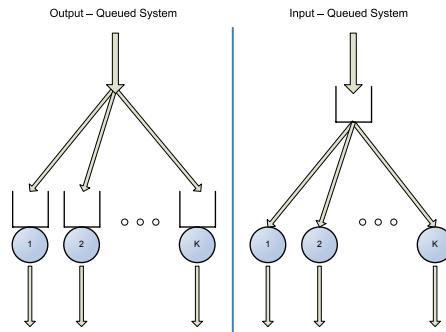


FIG 15. Output vs. Input-queued scheme

Output-queued schemes are inherently less efficient than their input-queued 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 \(2010\)](#); [Man-](#)

delbaum 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 many-server 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 highly 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? This should be tested empirically, as was done to a certain extent in Tseytlin and Zviran (2008).

2. *The role of information availability in routing and its influence on transfer delays:* An additional important aspect of routing schemes, which directly affects patient delays, is the availability of information in the system, 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. 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. It is interesting to investigate the effect of information availability on system performance and fairness towards patients and medical staff.

4.5. *Fairness in the ED-to-IW process.* Transfer policies may have ramifications on issues related to fairness towards *customers* (patients) and towards *servers* (medical and nursing staff). We investigate both aspects in the next several subsections.

4.5.1. *Fairness towards patients.* In Section 4.4 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 considerations, within each patient type there were significant FCFS violations as well. Specifically, 31% were overtaken by at least one patient of the same type, most of them not within the same ward, and hence these violations are likely due to the output-queued scheme.

TABLE 2
Percentage of FCFS violations per type within each IW

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%	

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.

4.5.2. 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 §3.3.1) or patients delay. For the latter, [Chan, Armony and Bambos \(2011\)](#) show that such fairness is achieved via Maximum Weighted Matching.

4.5.3. Fairness towards staff. In Section 4.4 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 already mentioned in §3.2.1, 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 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 a higher flux. Since nurses and doctors are assigned to particular wards, the load on 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 unfair towards servers.

Our data have 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. The authors in Tseytlin and Zviran (2008) perform a *simulation* study of the *output-queued* system under various routing schemes. They proposed an algorithm that balances a weighted function of occupancy and flux to achieve both fairness and short delays.

4.5.4. 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 across medical staff. A prerequisite for solving the “fairness problem” is then to define and calculate workload appropri-

ately. As argued in Section 6.2, such calculations must include not only direct time per resource but also emotional and cognitive efforts, as well as other relevant factors. For example, the mix of medical conditions and patient severity might 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), or the more severe patients who would challenge the medical staff, and provide them with further learning and research opportunities; the latter is especially relevant in teaching hospitals such as Rambam.

5. A system view. In Sections 2, 3, and 4 we treated the three network components (ED, IWs, transfers) unilaterally. In contrast, we now underscore the importance of looking at this network 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 needed here.

It is natural to expect that the higher the occupancy in the IWs the longer the delays in transfer, due to limited IW resources. The left diagram in Figure 16 displays the average delay in transfer 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, 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 the fact that IW physicians perform their morning rounds at this time and cannot admit new patients. Then 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 1–3pm. 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 3–4pm when most of the IWs discharges are complete. Shi et al. (2012), in their study of a Singapore hospital, discover similar discharge patterns, though the discharge procedure there occurs somewhat earlier in the day.

Further data analysis reveals that patients that are transferred to the IWs before 9am have significantly shorter LOS; early hospitalization reduces

ALOS by 1 day. Thus, we argue that it is extremely important to shorten the ED-to-IW transfer process and improve the admission process in the IWs so that the first day of hospitalization is not “wasted”.

In Section 4.3, we discussed how transfer delays 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.

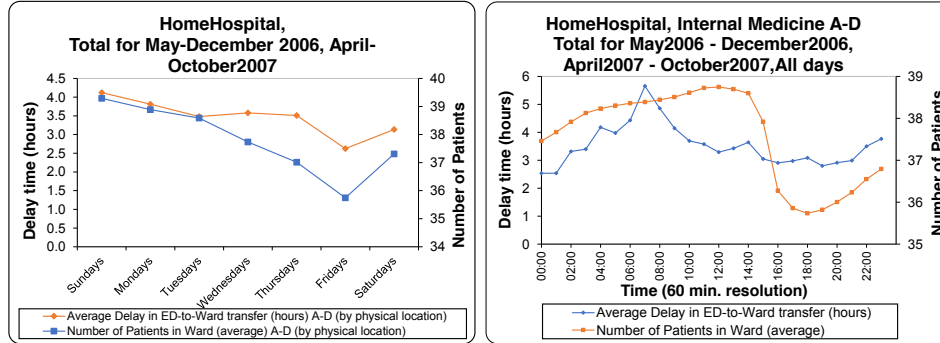


FIG 16. ED-to-IW transfer delays and number of patients in IW

5.1. Research opportunities. Our discussion suggests that daily routines (schedules) in the IWs have significant impact on transfer delays and thereby on ED workload. The question arises as to how one might wish to change these daily routines in view of this impact. The question fits well within a queuing context. The present daily routine at Rambam may be viewed as a priority scheme where current IW patients enjoy priority during morning physicians’ rounds, and discharged patients obtain priority in the afternoon, followed by newly-admitted patients. Is it possible to positively affect system performance by altering these priorities? More broadly, the challenge is to design priority schemes for a time-varying queuing network.

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 patients that are boarding in EDs (Figure 16) or in ICUs (Long and Mathews, 2012). Such boarding delays 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.

6. Discussion and concluding remarks. We have described research opportunities that arose from EDA of patient flow. We now highlight some common themes that have surfaced throughout the process. Specifically, we expand on the relationship between operational performance measures to overall hospital performance, discuss the concepts of workload and capacity and underline the importance of paying attention to time scales.

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 *clinical* but, practically, *operational* performance is the easiest to quantify, measure, track online and react upon. Moreover, operational performance is tightly coupled with the other dimensions, which explains its choice as a “language” that captures overall performance. For example, the fraction of patients who LWBS is a proxy for accessibility to care, and readmissions pertain to clinical quality of care.

Operational performance measures are often associated with patient flow. Among these, we discussed LWBS and “blocking” (Section 3.3.1), where patients end up being hospitalized in a ward different from that which is medically best for them (see Yom-Tov (2010) for more details); 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 require intra-ward data, which is beyond our 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, 2011) and the time-ingredients of an ED visit (treatment and waiting—for a resource, for synchronization or for a treatment to take effect; see Zaied (2011) and Atar, Mandelbaum and Zviran (2012)).

An operational measure that policy makers have been focusing on recently is readmissions. This is 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. It is commonly acknowledged, however, that one should consider readmissions judiciously as some of them

could be due to factors outside the hospital control (e.g. patients' own behavior, or care after discharge), or they may be an integral part of the treatment regiment. For example, returns within a few months to chemotherapy are typically *planned* and are unrelated to poor quality. But there are also 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 \(2011\)](#)) one should distinguish between these two readmission types by, for example, modeling planned (unplanned) readmissions as deterministic (stochastic) returns. Our hospital data supports the analysis of readmissions ([Mandelbaum et al., 2013](#)), which are further discussed in Section 4.2.2 of [EV](#). 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. We will return to time scales shortly.

6.2. Multi-dimensional workload. Operational performance of a service-system is determined by the gap, positive or negative, between its workload and the capacity assigned to process it. Workload is associated with a resource and, as such, used to gauge the requirements from that resource: for example, workload of nurses helps determine appropriate nurse staffing levels. Workload has also been shown to affect staff satisfaction level ([Aiken et al., 2002](#)) and patients' quality of care ([Batt and Terwiesch, 2012](#); [Kc and Terwiesch, 2009](#)).

In Queueing models, workload is typically an *average* quantity that is defined in steady-state: if λ is the arrival-rate of patients and S is the service time required from the nurse by a patient, then $R = \lambda \times E[S]$ is the workload of the nurse, which is commonly referred to as *offered load*. In *time-varying* environments ([Green, Kolesar and Whitt, 2007](#); [Reich, 2011](#)), notably hospitals, the offered-load $R(\cdot)$ is defined through the time-varying Little's law ([Bertsimas and Mourtzinou, 1997](#); [Green, Kolesar and Whitt, 2007](#)):

$$R(t) = \int_0^t \lambda(u) P\{S > t - u\} du, \quad t \geq 0,$$

in which $\lambda(\cdot)$ is the time-varying arrival rate (some mild assumptions are required for the integral to make sense).

As workload is matched against capacity, it must be measured in operational units. However, the workload of a nurse is affected by various factors beyond the mere time-content of nurse's tasks. For example, 1-minute of a standard chore does not compare with a 1-minute life-saving challenge. The

calculation of nurses' workload must therefore accommodate operational, emotional and cognitive factors, yet the outcome must be in standardized units that are "translatable" into staffing levels (Plonski et al., 2013).

It follows from the above that a comprehensive definition of personnel workload is inevitably complex. In the context of a medical ward, a natural approximation to workload is *relative occupancy*, namely the number of hospitalized patients in a ward, divided by the number of its beds. Assuming a constant beds-to-nurse ratio (Jennings and de Véricourt, 2011), homogeneous patient mix and an even distribution of workload over a patient's stay, occupancy is then well correlated with the *daily-routine* workload. However, as already discussed in Section 4.5.3, the level of medical attention that patients require typically declines during their stay. In addition, routine chores are less work intensive than *admissions and discharges* of patients, and the latter are naturally associated with patient *turnover* or *flux*, as opposed to bed occupancy. A proxy for workload must, therefore, acknowledge both occupancy and turnover, and possibly be also sensitive to the "age" (time since admission) of hospitalized patients.

Hence, we argue that the standard approach (presented above) for workload might be too simplified for the hospital environment. Here one must account for time-varying multi-dimensional aspects, which calls for research that aims at understanding, quantifying and measuring workload in health-care.

6.3. *Capacity.* Offered load characterizes operational demand for service which, in turn, must be matched by supply of the appropriate capacity. 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. This distinction has operational and cost accounting implications in 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 it is *flexible (controllable)*, both level- and hence cost-wise.

The association of flexible capacity with variable costs plays an important role in the Accounting/Strategic view of a hospital; Kaplan and Porter (2011) argue that most hospital costs are mistakenly judged as fixed while they ought to be viewed as variable costs, which means that the corresponding resource levels are flexible. This is an important observation, as it renders

controllable most resources in a hospital.

Our final point pertains to the characterization of capacity when it is flexible. Consider, for example, determining the capacity of the Ophthalmology ward in Rambam hospital, which happens to admit overflow patients from Internal wards (analogously to cross-trained agents in a call center; see [Aksin, Karaesmen and Ormeci \(2007\)](#)). Practically, this entails allocation of appropriate equipment, training of Ophthalmology ward personnel to be able to cater to IW patients and developing protocols for overflow of IW patients to Ophthalmology. Then the (dynamic) capacity of the Ophthalmology ward, if measured in patients-per-day say, would depend on its patient mix. The latter depends on the routing protocol which, in turn, determines the offered load. It follows that capacity and protocols better be determined jointly (unless their decoupling can be justified, as in [Armony, Gurvich and Mandelbaum \(2008\)](#)).

6.4. *Time-scales.* When analyzing ED-to-IWs flow (§4), the wards operate naturally on a time-scale of days while the ED time scale is hours. It follows that the wards serve as a random environment for the ED ([Ramakrishnan, Sier and Taylor, 2005](#)). Figure 9 (§3.2) reveals that the hourly scale is also of interest for IWs. These empirical examples arise from a service system that evolves in multiple time scales, which are all natural 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. \(2012\)](#), [Gurvich and Perry \(2012\)](#) and [Zacharias and Armony \(2013\)](#)).

The analysis—theoretical as well as empirical—of hospital units that operate under multiple time-scales offers further significant research opportunities. Hierarchical control of a hospital (e.g. strategic, tactical, and operational) is one such example, in which the strategic level generates constraints for tactical decisions which, in turn, percolate down to the operational level (e.g. [Zeltyn et al. \(2011\)](#)). To be specific, the number of beds in a ward is determined strategically which, tactically and operationally, sets staffing levels for nurses ([Jennings and de Véricourt, 2011](#)). This relates to our discussion in §3.2, where time scales helped determine what constitutes a server—a bed or a nurse.

6.5. *Some concluding comments on data-based research—a great opportunity but no less of a challenge.* Healthcare operations research is appropriately booming, and queueing applications to patient flow are trying to follow suit. In this context, the goal of the present work has been two-fold: first,

to encourage and strengthen, through data and its EDA, the natural link between theory and applications; and second, 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 (big) data is brought into the picture. Specifically, 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 inevitably time-consuming treatment. Next, using data in a single project, or a few for that matter, is wasteful—on the other hand, data-reuse and sharing, across student generations or research groups, requires infrastructure, documentation, maintenance and coordination. Finally, healthcare data is often confidential and proprietary, and that prevents reproducibility and slows down progress. (We will return to this last point momentarily.)

Fundamental changes are therefore essential—both within our OR/AP community as well as our potential healthcare partners: changes in education, organization and funding priorities, which takes us beyond our scope here. But we are 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 with providers of healthcare services, that are based on multidisciplinary (clinical, operational, financial, psychological) tracking of complete care-paths. Also, tracking resolution and scope will be at the level of the individual patient and provider, covering the full cycle of care.

6.5.1. *Towards a culture of reproducible research in empirical OR/AP.*

Data-based OR/AP research must strive for reproducibility of research outcomes—a fundamental principle in the traditional sciences. Reproducibility enables scrutiny of analysis and recommendations. This yields credibility and trust, which is an absolute prerequisite for influencing hospital practices.

Reproducible (Operations) Research is discussed in [Nestler \(2011\)](#), which is also a source for additional references and links. There have been some systematic attempts to establish a reproducibility culture in research ([Donoho et al., 2009](#)). It ought to start with funding agencies and journal policies: e.g. the Editorial Statement of the Finance Department in *Management Science* reads: “Authors of empirical and quantitative papers should provide or make

available enough information and data so that the results are reproducible.” It can advance with research such as [Karr \(2009\)](#), that aims at statistical analysis of distributed (unsharable) databases (e.g. hospital data); it will ideally culminate in a multitude of research labs, each providing free access to its data and serving its own research community and beyond.

A model for such a lab 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 [Nadjahrov et al. \(2013\)](#).

Acknowledgements. We dedicate our paper to the memory of the late David Sinreich. As a professor at the Technion IE&M, David was a pioneer and an example-to-follow, in acknowledging the necessity and importance of data- and simulation-based research in healthcare. In particular, the ED data collection that he orchestrated in 8 Israeli hospitals served as a proof-of-feasibility for our EDA.

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 (Director Rafi Beyar) and Technion IE&M Faculty (Former Dean Boaz Golany).

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

The authors, and the research presented here, owe a great deal to many additional individuals and institutions. This long list, which we can here acknowledge only in part, clearly must start with our host Rambam hospital—its capable management, dedicated nursing staff and physicians, and especially: Rafi Beyar, Rambam Director and CEO, who invited us at the outset to “use” the healthcare campus as our research laboratory; we gladly accepted the invitation, and have been since accompanied and advised by Zaher Azzam (Head of Rambam Internal Ward B), Sara Tzafrir (Head of Rambam Information Systems), and the OCR steering committee: Hana Adami (Head of Rambam Nursing), Yaron Barel (Rambam VP for Operations), Boaz Carmeli (IBM Research), Fuad Basis (Rambam ED), Danny Gopher (Technion IE&M), Avi Shtub (Technion IE&M), Segev Wasserkrug (IBM Research), Amir Weiman (Head of Rambam Accounting) and Pnina Vortman (IBM Research).

We are thankful to Technion students Kosta Elkin, Noga Rozenberg and Asaf Zviran: their projects and cooperation helped us analyze the ED-to-IWs process.

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 the

Statistics and Applied Mathematical Sciences Institute (SAMSI) of the NSF; the Department of Statistics and Operations Research (STOR), the University of North Carolina at Chapel Hill; the Department of Information, Operations and Management Sciences (IOMS), Leonard N. Stern School of Business, New York University; and the Department of Statistics, The Wharton School, University of Pennsylvania – the hospitality of all these institutions is acknowledged and 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. Last but certainly not least, we are grateful to the editor of *Stochastic Systems*, Peter Glynn, for leading and guiding us patiently and safely through the revision process.

REFERENCES

- ADLER, P. S., MANDELBAUM, A., NGUYEN, V. and SCHWERER, E. (1995). From Project to Process Management: An Empirically-Based Framework for Analyzing Product Development Time. *Management Science* **41** 458–484. 1.2
- AIKEN, L. H., CLARKE, S. P., SLOANE, D. M., SOCHALSKI, J. and SILBER, J. H. (2002). Hospital Nurse Staffing and Patient Mortality, Nurse Burnout, and Job Dissatisfaction. *JAMA* **288** 1987–1993. 6.2
- AKSIN, Z., ARMONY, M. and MEHROTRA, V. M. (2007). The Modern Call-Center: A Multi-Disciplinary Perspective on Operations Management Research. *Production and Operations Management* **16** 655–688. 1.2
- 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. 3.3.3, 6.3
- ALLON, G., DEO, S. and LIN, W. (2010). Impact of Size and Occupancy of Hospital on the Extent of Ambulance Diversion: Theory and Evidence. Working paper. 2.3
- ARMONY, M. (2005). Dynamic Routing in Large-Scale Service Systems with Heterogeneous Servers. *Queueing Systems* **51** 287–329. 1
- ARMONY, M., GURVICH, I. and MANDELBAUM, A. (2008). Service-Level Differentiation in Call Centers with Fully Flexible Servers. *Management Science* **54** 279–294. 6.3
- ARMONY, M. and WARD, A. (2010). Fair Dynamic Routing in Large-Scale Heterogeneous-Server Systems. *Operations Research* **58** 624–637. 4.5.3, 4.5.4
- ARMONY, M., ISRAELIT, S., MANDELBAUM, A., MARMOR, Y. N., TSEYTLIN, Y. and YOM-TOV, G. B. (2013). Patient Flow in Hospitals: A Data-Based Queueing-Science Perspective. An Extended Version (EV). 1.1, 1.3, 1.4, 2, 2.2.2, 2.2.2, 3.2.1, 3.3.1, 3.3.4, 4.1, 2, 6.1
- ATAR, R., MANDELBAUM, A. and ZVIRAN, A. (2012). Control of Fork-Join Networks in Heavy Traffic. Allerton Conference,. 4.2.1, 6.1
- ATAR, R. and SHWARTZ, A. (2008). Efficient Routing in Heavy Traffic under Partial Sampling of Service Times. *Mathematics of Operations Research* **33** 899–909. 1
- 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. 3.3.3
- BALASUBRAMANIAN, H., BANERJEE, R., DENTON, B., NAESSENS, J., WOOD, D. and STAHL, J. (2010). Improving clinical access and continuity using physician panel redesign. *Journal of General Internal Medicine* **25** 1109–1115. 3.3.3
- BATT, R. J. and TERWIESCH, C. (2012). Doctors Under Load: An Empirical Study of State-Dependent Service Times in Emergency Care. Working paper. 2.2.3, 6.2
- BEKKER, R. and DE BRUIN, A. M. (2010). Time-Dependent Analysis for Refused Admissions in Clinical Wards. *Annals of Operations Research* **178** 45–65. 3.2.1
- 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. 2.3
- BERTSIMAS, D. and MOURTZINO, G. (1997). Transient Laws of Non-stationary Queueing Systems and their Applications. *Queueing Systems* **25** 115–155. 2.2.1, 6.2
- 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. 1.5

- BRILLINGER, D. (2002). John Wilder Tukey (1915–2000). *Notices of the American Mathematical Society* 193–201. [1.2](#)
- 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. [1.2](#), [2.3](#), [3.2](#), [3.2.1](#)
- BURSTRÖM, L., NORDBERG, M., ORNUNG, G., CASTRÉN, M., WIKLUND, T., ENGSTRÖM, M. L. and ENLUND, M. (2012). Physician-Led Team Triage Based on Lean Principles May be Superior for Efficiency and Quality? A Comparison of Three Emergency Departments with Different Triage Models. *Scandinavian Journal Trauma, Resuscitation, and Emergency Medicine* **20** 20–57. [2.3](#)
- CAMERER, C. F., LOEWENSTEIN, G. and RABIN, M., eds. (2003). *Advances in Behavioral Economics*. Princeton University Press. [1.2](#)
- 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. [4.3](#)
- CHAN, C., ARMONY, M. and BAMBOS, N. (2011). Fairness in Overloaded Parallel Queues. Working paper. [4.5.2](#)
- CHAN, C., YOM-TOV, G. B. and ESCOBAR, G. (2011). When to use Speedup: An Examination of Service Systems with Returns. Working paper. [1.5](#)
- CHEN, C., JIA, Z. and VARAIYA, P. (2001). Causes and Cures of Highway Congestion. *Control Systems, IEEE* **21** 26–33. [2.2.3](#), [4.3](#)
- CHEN, H. and YAO, D. D. (2001). *Fundamentals of Queuing Networks: Performance, Asymptotics, and Optimization*. Springer. [3.3](#)
- CHEN, H., HARRISON, J. M., MANDELBAUM, A., VAN ACKERE, A. and WEIN, L. (1988). Empirical Evaluation of a Queueing Network Model for Semiconductor Wafer Fabrication. *Operations Research* **36** 202–216. [1.2](#)
- 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. [1.5](#)
- 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 Queueing Theory. *Health Care Management Science* **10** 125–137. [1.5](#)
- 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. [1.5](#), [2](#), [3.2.1](#), [3.3.1](#)
- DENTON, B. T., ed. (2013). *Handbook of Healthcare Operations Management: Methods and Applications*. Springer. [1.5](#)
- DOBSON, G., TEZCAN, T. and TILSON, V. (2013). Optimal Workflow Decisions for Investigators in Systems with Interruptions. *Management Science*. Forthcoming. [2.3](#)
- DONOHO, D. L., MALEKI, A., SHAHRAM, M., RAHMAN, I. U. and STODDEN, V. (2009). Reproducible Research in Computational Harmonic Analysis. *IEEE Computing in Science & Engineering* **11** 8–18. [6.5.1](#)
- ELKIN, K. and ROZENBERG, N. (2007). Patients Flow from the Emergency Department to the Internal Wards. IE&M project, Technion (In Hebrew). [3.1](#), [4.5.3](#)

- 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. [2.2.1](#)
- GANS, N., KOOLE, G. and MANDELBAUM, A. (2003). Telephone Call Centers: Tutorial, Review and Research Prospects. *Manufacturing, Services and Operations Management* **5** 79–141. [1.2](#), [3.3](#)
- GERLA, M. and KLEINROCK, L. (1980). Flow Control: A Comparative Survey. *IEEE Transactions on Communications* **28** 553–574. [2.2.3](#), [4.3](#)
- GITTINS, J., GLAZEBROOK, K. and WEBER, R. (2011). *Multi-armed Bandit Allocation Indices*. Wiley. [2.3](#)
- 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. [1.5](#)
- 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. [1.5](#)
- 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. [1.5](#), [2.2.1](#), [3.2.1](#), [6.2](#)
- GREEN, L. and YANKOVIC, N. (2011). Identifying Good Nursing Levels: A Queuing Approach. *Operations Research* **59** 942–955. [1.5](#), [3.2.1](#)
- 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. [2](#)
- GURVICH, I. and PERRY, O. (2012). Overflow Networks: Approximations and Implications to Call-Center Outsourcing. *Operations Research* **60** 996–1009. [6.4](#)
- GURVICH, I. and WHITT, W. (2010). Service-Level Differentiation in Many-Server Service Systems via Queue-Ratio Routing. *Operations Research* **58** 316–328. [1](#)
- 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. [2.3](#)
- HALL, R. W., ed. (2006). *Patient Flow: Reducing Delay in Healthcare Delivery*. Springer. [1.5](#), [2](#)
- HALL, R. W., ed. (2012). *Handbook of Healthcare System Scheduling*. Springer. [1.5](#)
- HALL, R., BELSON, D., MURALI, P. and DESSOUKY, M. (2006). Modeling Patient Flows Through the Healthcare System. In *Patient Flow: Reducing Delay in Healthcare Delivery* (R. W. Hall, ed.) 1 1–45. Springer. [1.5](#)
- HERMAN, R. (1992). Technology, Human Interaction, and Complexity: Reflections on Vehicular Traffic Science. *Operations Research* **40** 199–211. [1.2](#)
- 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. [2.3](#)
- HUANG, J. (2013). Patient Flow Management in Emergency Departments PhD thesis, National University of Singapore (NUS). [1.5.1](#)
- HUANG, J., CARMELI, B. and MANDELBAUM, A. (2011). Control of Patient Flow in Emergency Departments: Multiclass Queues with Feedback and Deadlines. Working paper. [1.5.1](#), [2.3](#), [6.1](#)

- 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. 2.3
- IHI, (2011). Patient First: Efficient Patient Flow Management Impact on the ED. *Institute for healthcare improvement*. <http://www.ihl.org/knowledge/Pages/ImprovementStories/PatientFirstEfficientPatientFlowManagementED.aspx>. 2
- 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. 3.3
- JCAHO, (2004). JCAHO Requirement: New Leadership Standard on Managing Patient Flow for Hospitals. *Joint Commission Perspectives* **24** 13–14. 1.1
- JENNINGS, O. B. and DE VÉRICOURT, F. (2008). Dimensioning Large-Scale Membership Services. *Operations Research* **56** 173–187. 1.5
- JENNINGS, O. B. and DE VÉRICOURT, F. (2011). Nurse Staffing in Medical Units: A Queueing Perspective. *Operations Research* **59** 1320–1331. 1.5, 3.2.1, 3.3, 6.2, 6.4
- 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. 3.3.3
- KAPLAN, R. S. and PORTER, M. E. (2011). How to Solve the Cost Crisis in Health Care. *Harvard Business Review* **89** 46–64. 6.3
- KARR, A. F. (2009). Secure Statistical Analysis of Distributed Databases, Emphasizing What We Don't Know. *Journal of Privacy and Confidentiality* **1**. <http://repository.cmu.edu/jpc/vol1/iss2/5>. 6.5.1
- 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. 1.5, 2.2.3, 6.2
- KELLY, F. P. (1979). *Markov Processes and Reversibility*. Wiley. 2.2.2
- KIM, S. H. and WHITT, W. (2013). Are Call Center and Hospital Arrivals Well Modeled by Nonhomogeneous Poisson Processes? Working paper, Columbia University. 2.2.1
- LEITE, S. C. and FRAGOSO, M. D. (2013). Diffusion Approximation for Signaling Stochastic Networks. *Stochastic Processes and their Applications* **123** 2957–2982. 4.2.1
- LELAND, W. E., TAQQU, M. S., WILLINGER, W. and WILSON, D. V. (1994). On the Self-Similar Nature of Ethernet Traffic (Extended Version). *Operations Research* **2** 1–15. 1.2
- LONG, E. F. and MATHEWS, K. M. (2012). “Patients Without Patience”: A Priority Queueing Simulation Model of the Intensive Care Unit. Working paper. 4.3, 5.1
- MAA, J. (2011). The Waits that Matter. *The New England Journal of Medicine*. 4.3
- 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. 1.5.1, 2.2.1
- MAMAN, S., ZELTYN, S. and MANDELBAUM, A. (2011). Uncertainty in the Demand for Service: The Case of Call Centers and Emergency Departments. Working paper. 1.5.1
- 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. 1.5.1, 3.2.1, 4.5.1, 4.5.3, 4.5.4, 6.4

- MANDELBAUM, A., SAKOV, A. and ZELTYN, S. (2000). Empirical Analysis of a Call Center. Technical Report, <http://iew3.technion.ac.il/serveng/References/ccdata.pdf>. 1.2
- 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. 1
- 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. 6.1
- MARMOR, Y. N. (2003). Developing a Simulation Tool for Analyzing Emergency Department Performance. Master’s thesis, Technion—Israel Institute of Technology. 1.3
- 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. 1.5.1
- MARMOR, Y. N., ROHLEDER, T., HUSCHKA, T., COOK, D. and THOMPSON, J. (2011). Cardio Vascular Surgery Simulation in Support of System Engineering Decision Making. In preparation. 2.2.3
- 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. 1.5.1, 2.3
- McHUGH, M., VAN DYKE, K., McCLELLAND, M. and MOSS, D. (2011). Improving Patient Flow and Reducing Emergency Department Crowding. *Agency for health-care research and quality*. <http://www.ahrq.gov/research/findings/final-reports/ptflow/index.html>. 2
- MEDICARE USA, (2013). Hospital Compare: 30-Day Death and Readmission Measures Data. <http://www.medicare.gov/HospitalCompare/Data/RCD/30-day-measures.aspx>. 6.1
- NADJHAHROV, E., TROFIMOV, V., GAVAKO, I. and MANDELBAUM, A. (2013). Home-Hospital (Rambam): EDA via SEESStat 3.0 to Reproduce “On Patients Flow in Hospitals”. http://ie.technion.ac.il/Labs/Serveng/files/HHD/reproducing_flow_paper.pdf. 6.5.1, 6.5.1
- NESTLER, S. (2011). Reproducible (Operations) Research: A Primer on Reproducible Research and Why the O.R. Community Should Care About it. *ORMS-Today* **38**. 6.5.1
- NGUYEN, V. (1994). The Trouble with Diversity: Fork-Join Networks with Heterogeneous Customer Population. *The Annals of Applied Probability* 1–25. 4.2.1
- OREDSSON, S., JONSSON, H., ROGNES, J., LIND, L., GÖRANSSON, K. E., EHRENBORG, A., ASPLUND, K., CASTRÉN, M. and FARROHKHIA, N. (2011). A Systematic Review of Triage-Related Interventions to Improve Patient Flow in Emergency Departments. *Scandinavian Journal Trauma, Resuscitation, and Emergency Medicine* **July 19** 19–43. 2.3
- 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. 3.2.1, 6.2
- 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. 6.4
- RAMBAM, Rambam Health Care Campus, Haifa, Israel. <http://www.rambam.org.il/Home+Page/>. 1.4

- RAMBAMDATA, Rambam Hospital data repositories. Technion SEELab, <http://seeserver.iem.technion.ac.il/databases/>. 6.5.1, 6.5.1
- REICH, M. (2011). The Workload Process: Modelling, Inference and Applications. Master's thesis, Technion—Israel Institute of Technology. 6.2
- SEELAB, SEE Lab, Technion—Israel Institute of Technology. <http://ie.technion.ac.il/Labs/Serveng/>. 1, 1, 6.5.1, 6.5.1
- SEESERVER, Server of the Center for Service Enterprise Engineering. <http://seeserver.iem.technion.ac.il/see-terminal/>. 6.5.1
- SEESTAT, SEESStat documentation, Technion—Israel Institute of Technology. <http://ie.technion.ac.il/Labs/Serveng/>. 1, 6.5.1, 6.5.1
- SHI, P., CHOU, M. C., DAI, J. G., DING, D. and SIM, J. (2012). Hospital Inpatient Operations: Mathematical Models and Managerial Insights. Working paper. 1.3, 1.5, 3.2, 3.2.1, 5, 6.4
- STOLYAR, S. (2005). Optimal Routing in Output-Queued Flexible Server Systems. *Probability in the Engineering and Informational Sciences* **19** 141–189. 1, 1
- 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. 2.2.3
- SUN, J. (2006). *The Statistical Analysis of Interval-Censored Failure Time Data*. Springer. 2.3
- TEZCAN, T. (2008). Optimal Control of Distributed Parallel Server Systems under the Halfin and Whitt Regime. *Math of Operations Research* **33** 51–90. 1
- THE OCR PROJECT (IBM-RAMBAM-TECHNION), (2011). Service Science in Hospitals: A Research-Based Partnership for Innovating and Transforming Patients Care. http://ie.technion.ac.il/Labs/Serveng/files/Project_summary_for_SRII.pdf. 1.2
- 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. 4.5.2
- THORIN, O. (1977). On the Infinite Divisibility of the Lognormal Distribution. *Scandinavian Actuarial Journal* **1977** 121–148. 3.2.1
- TSEYTLIN, Y. (2009). Queueing Systems with Heterogeneous Servers: On Fair Routing of Patients in Emergency Departments. Master's thesis, Technion—Israel Institute of Technology. 1.3, 1.5.1
- 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. 1, 4.5.3, 4.5.4
- TUKEY, J. W. (1977). *Exploratory Data Analysis*. Addison Wesley. 1.2
- WARD, A. and ARMONY, M. (2013). Blind Fair Routing in Large-Scale Service Systems with Heterogeneous Customers and Servers. *Operations Research* **61** 228–243. 4.5.4
- WHITT, W. (2012). Fitting Birth-and-Death Queueing Models to Data. *Statistics and Probability Letters* **82** 998–1004. 2.2.2
- YOM-TOV, G. B. (2010). Queues in Hospitals: Queueing Networks with ReEntering Customers in the QED Regime. PhD thesis, Technion—Israel Institute of Technology. 1.5.1, 3.2.1, 6.1
- YOM-TOV, G. B. and MANDELBAUM, A. (2011). Erlang-R: A Time-Varying Queue with ReEntrant Customers, in Support of Healthcare Staffing. Forthcoming MS&OM. 1.5, 1.5.1, 2.2.2, 2.3, 3.3, 6.1

- ZACHARIAS, C. and ARMONY, M. (2013). Joint Panel Sizing and Appointment Scheduling in Outpatient Care. Working paper. [6.4](#)
- 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. [5.1](#), [6.1](#)
- ZALENSKI, R. J., RYDMAN, R. J., TING, S., KAMPE, L. and SELKER, H. P. (1998). A National Survey of Emergency Department Chest Pain Centers in the United States. *The American Journal of Cardiology* **81** 1305–1309. [2.3](#)
- 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**. [1.5.1](#), [2.3](#), [6.1](#), [6.4](#)
- ZHU, B., ARMONY, M. and CHAN, C. W. (2013). Critical Care in Hospitals: When to Introduce a Step Down Unit? Working paper, Columbia University. [4.2](#)

Appendix: Accessing Data repositories and EDA tools at the SEELab. [SEELab](#) is a data-based research laboratory, residing at the IE&M Faculty of the Technion in Haifa, Israel. (SEE stands for “Service Enterprise Engineering”.) SEELab maintains a repository for 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 call-by-call multi-year data from 4 call centers, an internet academic website, 8 emergency departments (mainly their arrivals data) and 4 years of data from the Rambam Hospital—this is the empirical foundation for the present paper.

The EDA environment of SEELab is [SEESStat](#)—a software platform that enables real-time statistical analysis of service data at seconds-to-months time resolutions. SEESStat 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. SEESStat 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](#)).

Three SEELab data-bases are publicly accessible at the SEELab server [SEEServer](#): two from call centers and one from the Rambam hospital. For example, data from a U.S. banking call center covers the operational history of close to 220 million calls, over close to 3 years; 40 million of these calls were served by (up to 1000) agents and the rest by a VRU (answering machine). The Rambam data is described in [§1.4](#).

SEESStat Online: The connection protocol to SEELab data, for any re-

search 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 **SEESat 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 SEESat, SEELab's EDA environment, and via SEESat 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 SEESat and then let the tutorial guide you, hands-on, through SEESat basics—this should take no more than 1.5 hours.

On data cleaning and maintenance. There were plenty of records that were flawed due to archiving or simply system errors. These were identified via their inconsistency with trustable data and hence corrected or removed. But more challenging was the identification of records that had been included in the data due to some regulations, rather than physical transactions. For example, some unreasonable workload profiles led to the discovery of a high fraction of “transfers” from the ED to a virtual ward, all occurring precisely at 11:59pm; subsequent analysis managed to associate each of these transfers with a physical transfer, from the ED to some actual ward on the *following* day. The reason for the inclusion of such virtual transfers was financial, having to do with regulations of insurance reimbursement. And this is just one out of many examples.

Reproducing our EDA and beyond. Rambam data is publicly available, either for downloading ([RambamData](#) consists of records per individual customers) or through SEESat. The download link includes data documentation. To facilitate reproducibility, the document [Nadjahrov 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.

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

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

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

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

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

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