

Service Engineering in Action: The Palm/Erlang-A Queue, with Applications to Call Centers

Avishai Mandelbaum and Sergey Zeltyn

Faculty of Industrial Engineering & Management
Technion, Haifa 32000, ISRAEL

emails: avim@tx.technion.ac.il, zeltyn@ie.technion.ac.il

March 16, 2005

Abstract

Our note¹ is dedicated to the Palm/Erlang-A Queue. This is the simplest practice-worthy queueing model, that accounts for customers' impatience while waiting. The model is gaining importance in support of the staffing of call centers, which is a central step in their Service-Engineering. We discuss computations of performance measures, both theoretical and software-based (via the 4CallCenter software). Then several examples of Palm/Erlang-A applications are presented, mostly motivated by and based on real call center data.

Acknowledgements. The research of both authors was supported by ISF (Israeli Science Foundation) grants 388/99, 126/02 and 1046/04, by the Niderzaksen Fund and by the Technion funds for the promotion of research and sponsored research.

¹Parts of the text are adapted from [8], [15], [17] and [22]

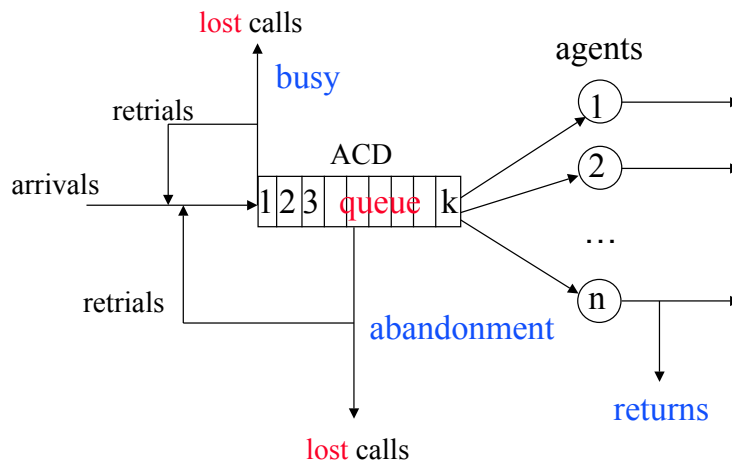
Contents

1	Introduction	1
2	Significance of abandonment in modelling and practice	2
3	Birth-and-death process representation	5
4	Operational measures of performance	6
4.1	Practical measures: accounting for Abandonment	6
4.2	Calculations: the 4CallCenters software	7
4.3	A general approach for computing operational performance measures	8
4.4	Relation between average wait and the fraction abandoning	9
5	Parameter estimation in a call center environment	10
6	Approximations	12
7	Applications to call centers	15
7.1	Erlang-A performance measures: comparison against real data	15
7.2	Erlang-A approximations: comparison against real data	16
8	Some advanced features of 4CallCenters	16
9	Some open research topics	19
9.1	Dimensioning the Erlang-A queue	19
9.2	Human behavior	21
9.3	Uncertainty in parameter values	23
A	The Erlang-A queue: useful formulae for the steady-state distribution and some performance measures	26

1 Introduction

Service Engineering is a newly emerging discipline that seeks to develop scientifically-based engineering principles and tools, often culminating in software, which support the design and management of service operations. **Contact Centers** are service organizations for customers who seek service via the phone, fax, e-mail, chat or other tele-communication channels. A particularly important type of contact centers are the **Call Centers**, which predominantly serve phone calls. Due to advances in Information and Communication Technology, the number, size and scope of contact centers, as well as the number of people who are employed there or use them as customers, grows explosively. For example, in the U.S. alone, the call center industry is estimated to employ several million agents which, in fact, outnumbers agriculture. In Europe, the number of call center employees in 1999-2000 was estimated, for example, by 600,000 in the UK (2.3% of the total workforce) and 200,000 in Holland (almost 3%) [3]. Bittner et al. [5] assess that, in Germany in 2001, there were between 300,000 to 400,000 (1-2%) employed in the call center industry.

Figure 1: Schematic representation of a telephone call center



In a large performance-leader call center, many hundreds of agents could serve many thousands of calling customers per hour; agents' utilization levels exceed 90%, yet about 50% of the customers are answered immediately upon calling; callers who are delayed demand a response within seconds, the vast majority gets it, and scarcely few of the rest, say 1% of those calling, abandon during peak-congestion due to impatience. But most call centers are far from achieving such levels of performance. To these, scientific models are prerequisites for climbing the performance ladder, and the model described in this paper, namely **Palm/Erlang-A**,

should constitute their starting point. See Gans, Koole and Mandelbaum [13] or Helber and Stolletz [18] for reviews of state-of-the-art of research on telephone call centers. In addition, Mandelbaum [20] provides a comprehensive bibliography, namely references plus abstracts, of call-center-related research papers.

Modelling a Call Center. A simplified representation of traffic flows in a call center is given in Figure 1. Incoming calls form a single queue, waiting for service from one of n statistically identical agents. There are $k + n$ telephone trunk-lines. These are connected to an Automatic Call Distributor (ACD) which manages the queue, connects customers to available agents, and also archives operational data. Customers arriving when all lines are occupied encounter a busy signal. Such customers might try again later (“retrial”) or give up (“lost call”). Customers who succeed in getting through at a time when all agents are busy (that is, when there are at least n but fewer than $k + n$ customers within the call center), are placed in the queue. If waiting customers run out of patience before their service begins, they hang up (“abandon”). After abandoning, customers might try calling again later while others are lost. After service, there are “positive” returns of satisfied customers, or “negative” returns due to complaints.

Note that the model in Figure 1 ignores multiple service types and skilled-based routing that are present in many modern call centers. However, a lot of interesting questions still remain open (see Section 9) even for models with homogeneous servers/customers.

In basic models, the already simple representation in Figure 1 is simplified even further. Specifically, in the present paper we assume that there are enough trunk-lines to avoid busy signals ($k = \infty$). This assumption prevails in today’s call centers. In addition, we assume out retrials and return calls, which corresponds to absorbing them within the arrivals. (See, for example, Aguir et al. [2] for an analysis that takes retrials into account.) However, and unlike most models used in practice, here we do acknowledge and accommodate abandonment. The reasons for this will become clear momentarily.

2 Significance of abandonment in modelling and practice

The classical M/M/ n queueing model, also called the **Erlang-C** model, is the one most frequently used in *workforce management* of call centers. Erlang-C assumes Poisson arrivals at a constant rate λ , exponentially distributed service times with a rate μ , and n independent statistically-identical agents. (Time-varying arrival rates are accommodated via piecewise constant approximations.) But Erlang-C does not allow abandonment. This, as will now be argued, is a significant deficiency: customer abandonment is not a minor, let alone a negligible, aspect of call center operations. We now support this last statement, first qualitatively and then quantitatively.

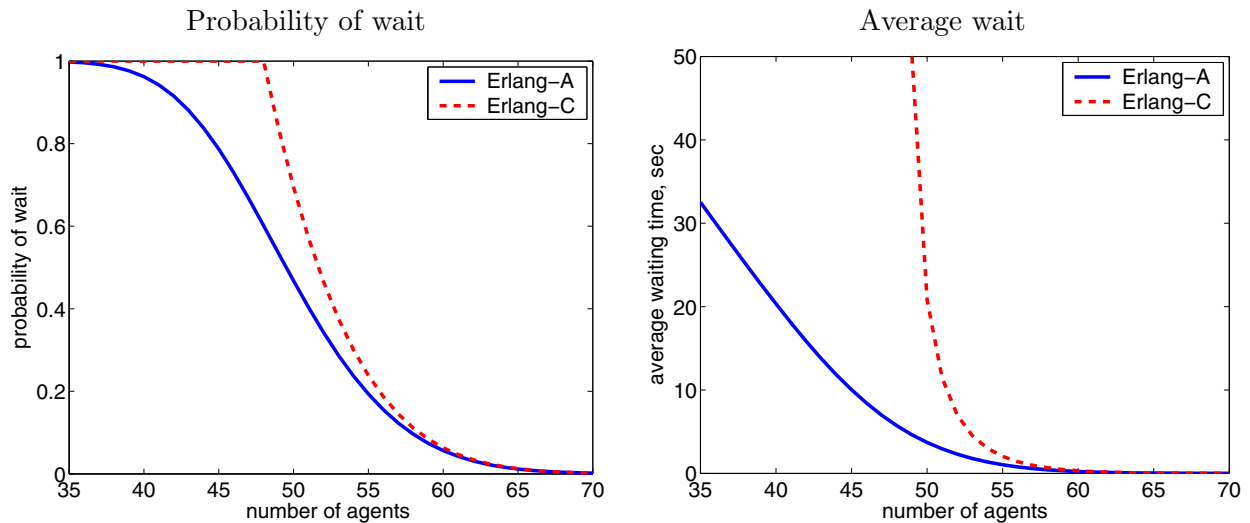
- Abandonment statistics constitute the only ACD measurement that is customer-subjective: those who abandon declare that the service offered is not worth its wait. (Other ACD data, such as average waiting times, are “objective”; they also do not include the only other customer-subjective operational measures, namely retrieval/return statistics.)
- Some call centers focus on the average waits of only those who get served, which does not acknowledge abandoning customers. But under such circumstances, the service-order that optimizes performance is LIFO = Last-In-First-Out [14], which clearly suggests that a distorted focus has been chosen.
- Ignoring abandonment can cause either under- or over-staffing: On the one hand, if service level is measured only for those customers who reach service, the result is unjustly optimistic - the effect of an abandonment is less delay for those further back in line, as well as for future arrivals. This would lead to under-staffing. On the other hand, using workforce management tools that ignore abandonment would result in over-staffing as actually fewer agents are needed in order to meet most abandonment-ignorant service goals.

The Palm/Erlang-A model: Palm [26] introduced a simple (tractable) way to model abandonment. He suggested to enrich Erlang-C (M/M/n) in the following manner. Associated with each arriving caller there is an exponentially distributed *patience time* with mean θ^{-1} . An arriving customer encounters an *offered waiting time*, which is defined as the time that this customer would have to wait given that her or his patience is infinite. If the offered wait exceeds the customer’s patience time, the call is then abandoned, otherwise the customer awaits service. The patience parameter θ will be referred to as the individual *abandonment rate*. (We shall omit “individual”, when obvious.) We denote this model by M/M/n+M, and refer to it as Palm/Erlang-A, or Erlang-A for short. Here the **A** stands for **A**bandonment, as well as for the fact that the model interpolates between Erlang-C and Erlang-B. (The latter is the M/M/n/n model, in which there are n trunk lines ($k=0$), hence customers that cannot be served immediately are blocked.)

With Erlang-A, the quantitative significance of abandonment can be demonstrated through simple numerical examples. We start with Figure 2, which shows the fraction of delayed customers and the average wait, when calculated via Erlang-C (M/M/n), and a corresponding Erlang-A (M/M/n+M) model. In both models, the arrival rate is 48 calls per minutes, the average service time equals 1 minute, and the number of agents is varied from 35 to 70. Average patience is taken to be 2 minutes for the Erlang-A model. Clearly, the two curves convey rather different pictures of what is happening in the system they depict, especially within the range of 40 to 50 agents: in particular, and as shown below, Erlang-C is stable only with 49 or more agents, while Erlang-A is always stable.

The above M/M/n and M/M/n+M models are further compared in Table 1. Note that exponential patience with an average of 2 minutes gives rise to 3.1% abandonment. Then note that the average wait and queue length are both strikingly shorter with only 3.1% abandonment taking place. Indeed, “The fittest survive” and wait less - much less. (Significantly, this high-level performance is *not* achieved if the arrival rate to the M/M/n system is decreased by 3.1%; for example, the “average speed of answer” in such a case is 8.8 seconds, compared with 3.7 seconds. The reason is that abandonment reduce workload precisely when needed, namely when congestion is high.) Finally, note that system performance in such heavy traffic is very sensitive to staffing levels. In our example, adding 3 or 4 agents (from 50 to say 54) to M/M/n would result in M/M/n+M performance, as emerging from the horizontal distance between the graphs in Figure 2. Nonetheless, since personnel costs are the major operational costs of running call centers (prevalent estimates run at about 60-75% of the total), even a 6%-8% reduction in personnel is economically significant (and much more so for large call centers that employ thousands of agents).

Figure 2: **Comparison between Erlang-A and Erlang-C** [15, 34]
 48 calls per min., 1 min. average service time, 2 min. average patience



As a final demonstration of the significance of abandonment, we now use it to explain a phenomenon that has puzzled queueing theorists: It is the observation that, in practice, simple deterministic approaches often lead to surprisingly good results. For example, consider a call center with averages of 6000 calls per hour and service time of 4 minutes. Such a call center gets an average of $(6000 : 60) \cdot 4 = 400$ minutes of work per minute. The deterministic approach then prescribes 400 service agents to cope with this load (1 agent-minute per 1 work-minute), which is a questionable recommendation according to standard queueing models. For example, Erlang-C

Table 1: **Comparing models with/without abandonment**
 50 agents, 48 calls per min., 1 min. average service time, 2 min. average patience

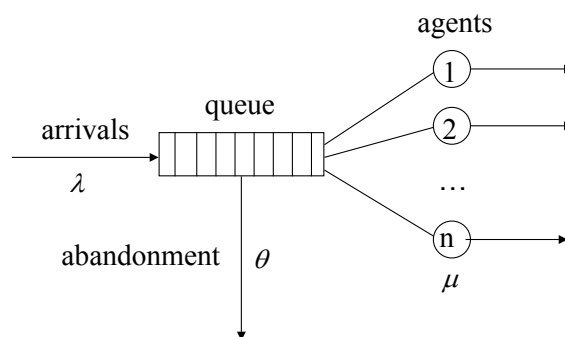
	M/M/n	M/M/n+M	M/M/n, $\lambda \downarrow 3.1\%$
Fraction abandoning	–	3.1%	-
Average waiting time	20.8 sec	3.7 sec	8.8 sec
Waiting time's 90-th percentile	58.1 sec	12.5 sec	28.2 sec
Average queue length	17	3	7
Agents' utilization	96%	93%	93%

would then be unstable, and its waiting times and queue-lengths would increase indefinitely. But now assume that customers abandon, as they actually do, and assign a reasonable parameter to their average patience, say equal to the average service time. Then, under Erlang-A, about 50% of the customers would be answered immediately upon calling, the average wait would be a mere 5 seconds, agents' utilization would be 98%, and all this at the cost of 2% abandonment – a remarkable performance indeed. (See the Remark in Section 6 for a more formal explanation.)

3 Birth-and-death process representation

Figure 3 provides a representation of the traffic flows in Erlang-A, and a comparison with Figure 1 clearly reveals its limitations. (Nevertheless, and as we hope to demonstrate, Erlang-A still turns out very useful and insightful, both theoretically and practically.)

Figure 3: **Schematic representation of the Erlang-A model**



Erlang-A is characterized by 4 parameters: λ , μ , θ and n . Here λ is the calling rate (calls per unit of time); μ is the service rate ($1/\mu$ is the average duration of service); $1/\theta$ is the average patience of a customer; and n is the number of servers/agents. More formally, in the Erlang-A model customers arrive to the queueing system according to a Poisson(λ) process. Customers are

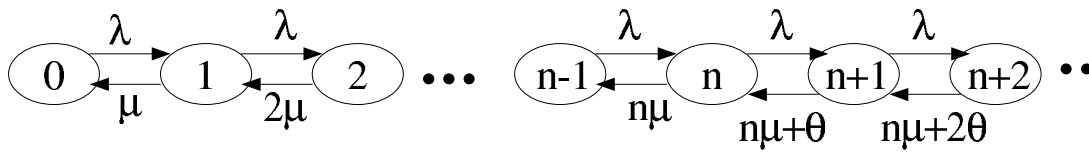
equipped with *patience times* τ that are $\exp(\theta)$, i.i.d. across customers. And service times are i.i.d. $\exp(\mu)$. Finally, the processes of arrivals, patience and service are mutually independent.

For a given customer, the patience time τ is the time that the customer is *willing* to wait for service - a wait that reaches τ results in an abandonment. Let V denote the *offered waiting time* - the time a customer, equipped with infinite patience, *must* wait in order to get service. The actual waiting/queueing time then equals

$$W = \min\{V, \tau\}.$$

Denote by $L(t)$ the number-in-system at time t (includes both customers being served and waiting in the queue). Then $L = \{L(t), t \geq 0\}$ is a Markov birth-and-death process, with the following transition-rate diagram:

Figure 4: **Transition-rate diagram of the Erlang-A model**



An analysis of a birth-and-death process usually starts with verifying that it reaches steady-state (it always does, in our case), and it then continues with calculation of its *limiting/steady-state distribution*, defined by:

$$\pi_j \triangleq \lim_{t \rightarrow \infty} P\{L(t) = j\}, \quad j = 0, 1, 2, \dots \quad (3.1)$$

Alternatively, π_j can be characterized as the fraction of time that the system spends in state j , when in steady-state. Formulae for the steady-state distribution of Erlang-A are presented in the Appendix.

4 Operational measures of performance

In order to understand and apply the Erlang-A model, one must first define its measures of performance, and then be able to calculate them. Moreover, since a call center can get very large (thousands of agents), the implementation of these calculations must be both fast and numerically stable.

4.1 Practical measures: accounting for Abandonment

The most popular measure of operational (positive) performance is the fraction of served customers that have been waiting less than some given time, or formally $P\{W \leq T, \text{Sr}\}$, where W

is the (random) waiting time in steady-state, $\{Sr\}$ is the event “customer gets service” and T is a target time that is determined by Management/Marketing. However, as explained before, performance measures must take into account those customers who abandon. Indeed, if forced into choosing a *single* number as a proxy for operational performance, we recommend the probability to abandon $\mathbf{P}\{\mathbf{Ab}\}$, the fraction of customers who explicitly declare that the service offered is not worth its wait. Some managers actually opt for a refinement that excludes those who abandon within a very short time, formally $P\{W > \epsilon; Ab\}$, for some small $\epsilon > 0$, e.g. $\epsilon = 3$ seconds. The justification is that those who abandon within 3 seconds can not be characterized as poorly served. There is also a practical rationale that arises from physical limitations, specifically that such “immediate” abandonment could in fact be a malfunction or an inaccuracy of the measurement devices.

The single abandonment measure $P\{Ab\}$ can be in fact refined to account explicitly for those customers who were or were not well-served. To this end, we propose the following four-dimensional service measure, given 2 parameters T and ϵ :

- $P\{W \leq T; Sr\}$ - fraction of well-served;
- $P\{W > T; Sr\}$ - fraction of served, with a potential for improvement;
- $P\{W > \epsilon; Ab\}$ - fraction of poorly-served;
- $P\{W \leq \epsilon; Ab\}$ - fraction of those whose service-level is undetermined - see the above for an elaboration.

Our proposed 4-component measure is not commonly used and most workforce management software tools are incapable of calculating it. To have it practical, we now describe how it can be implemented via the software tool 4CallCenters [12].

4.2 Calculations: the 4CallCenters software

Black-box Erlang-A calculations, as well as many other useful features, are provided by the free-to-use software 4CallCenters [12]. (This software is being regularly upgraded.) The calculation methods are described in Appendix B of [15]; they were developed in the Technion’s M.Sc. thesis of the first author, Ofer Garnett.

Figure 5 displays a 4CallCenters output and demonstrates how to calculate the four-dimensional service measure, introduced in Subsection 4.1.

The values of the four Erlang-A parameters are displayed in the middle of the upper half of the screen. Let $T = 30$ seconds and $\epsilon = 10$ seconds. Then one should perform computations twice: with *Target Time* 30 and 10 seconds. (Both computations appear in Figure 5.) We get:

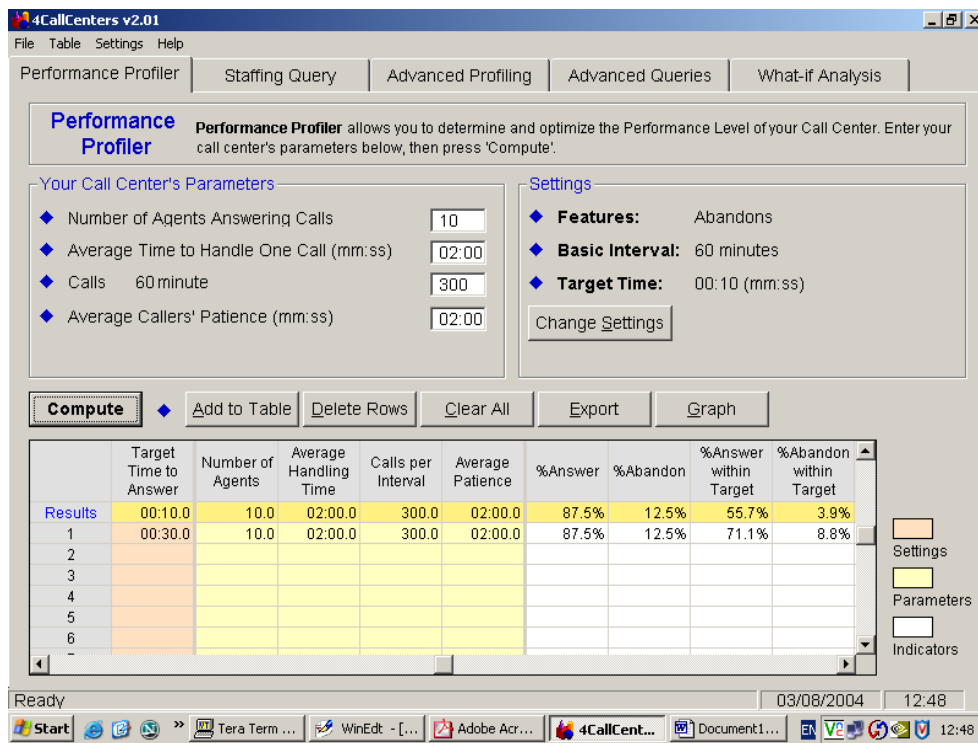
- $P\{W \leq T; Sr\}$ - fraction of well-served is equal to 71.1%;

- $P\{W > T; Sr\}$ - fraction of served, with a potential for improvement, is 16.4% (87.5% – 71.1%);
- $P\{W > \epsilon; Ab\}$ - fraction of poorly-served is 8.6% (12.5% – 3.9%);
- $P\{W \leq \epsilon; Ab\}$ - fraction of those whose service-level is undetermined is 3.9%.

Note that the 4CallCenters output includes many more performance measures than those displayed in Figure 5: one could scroll the screen to values of agents' occupancy, average waiting time, average queue length, etc.

In Section 8 we describe several examples of the more advanced capabilities of 4CallCenters.

Figure 5: 4Callcenters. Example of output.



4.3 A general approach for computing operational performance measures

Some explicit expressions of Erlang-A performance measures are provided in the Appendix. (See also Riordan [27].) However, we recommend to use more general M/M/n+G formulae, as the main alternative to 4CallCenters software. Indeed, Erlang-A is a special case of the M/M/n+G queue, in which patience times are generally distributed. A comprehensive list of M/M/n+G

formulae, as well as guidance for their application, appears in Mandelbaum and Zeltyn [24]. The preparation of [24] was triggered by a request from a large U.S. bank. Consequently, this bank has been routinely applying Erlang-A in the workforce management of its 10,000 telephone agents, who handle close to 150 millions calls yearly.

The handout [24] also explains how to adapt the M/M/n+G formulae to Erlang-A, in which patience is exponentially distributed. Specifically, see Sections 1,2 and 5 of [24].

4.4 Relation between average wait and the fraction abandoning

A remarkable property of Erlang-A, which in fact generalizes to other models with patience that is $exp(\theta)$, is the following linear relation between the fraction abandoning $P\{Ab\}$ and average wait $E[W]$:

$$P\{Ab\} = \theta \cdot E[W]. \tag{4.1}$$

Proof: The proof is based on the balance equation

$$\theta \cdot E[Q] = \lambda \cdot P\{Ab\}, \tag{4.2}$$

and on Little's formula

$$E[Q] = \lambda \cdot E[W], \tag{4.3}$$

where Q is the steady-state queue length. The balance equation (4.2) is a steady-state equality between the rate that customers abandon the queue (left hand side) and the rate that abandoning customers (i.e. - customers who eventually abandon) enter the system. Substituting Little's formula (4.3) into (4.2) yields formula (4.1). ■

Figure 6: Probability to abandon vs. average waiting time [8]

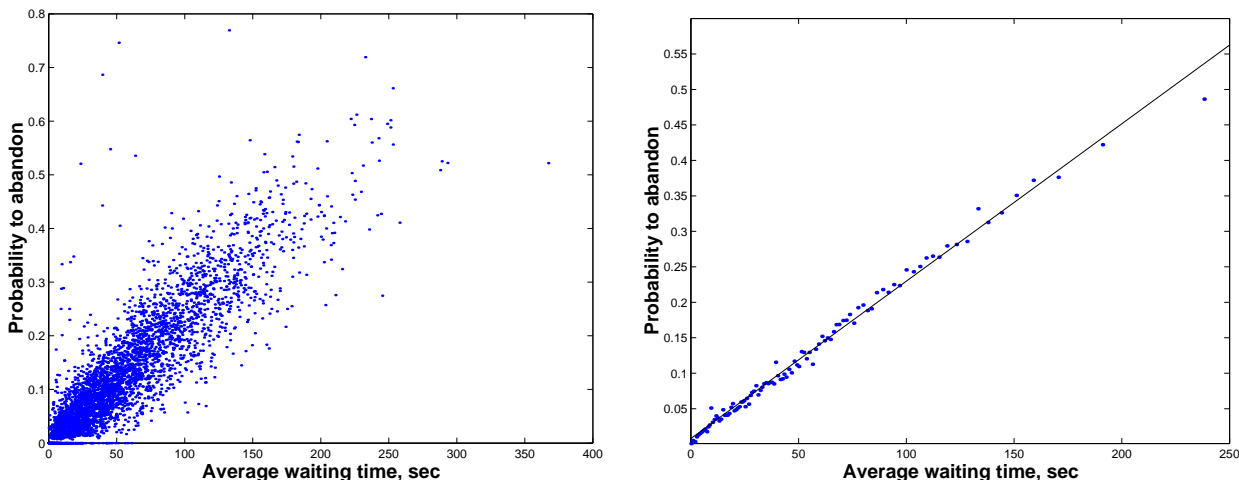


Figure 6 illustrates the relation (4.1). It was plotted using yearly data of an Israeli bank call center [9], which is analyzed in a Service Engineering course that is taught at the Technion [28, 9]. (See also Brown et al. [8] for statistical analysis of this call center data.) First, $P\{Ab\}$ and $E[W]$ were computed for the 4158 hour intervals that constitute the year 1999. The left plot of Figure 6 presents the resulting “cloud” of points, as they scatter on the plane. For the right plot, we are using an aggregation procedure that is designed to emphasize dominating patterns. Specifically, the 4158 intervals were ordered according to their average waiting times, and adjacent groups of 40 points were aggregated (further averaged): this forms the 104 points of the second plot in Figure 6. (The last point of the aggregated plot is an average of only 38 hour intervals.)

We observe a convincing linear relation (line) between $P\{Ab\}$ and $E[W]$. Based on (4.1) and Figure 6, the slope of this line is an estimate of the average patience, which here equals 446 seconds. In Brown et al. [8] and Mandelbaum et al. [19] it is shown that some $M/M/n+M$ assumptions do not prevail for the data [9]. Although arrivals are essentially Poisson, the service times are not exponential (in fact, they are very close to being lognormal). Patience times were shown to be non-exponential either. Yet, Erlang-A is proved useful for the performance analysis, which we demonstrate further in Section 7.

It is therefore important to understand the circumstances under which one can practically use simple relations that, theoretically, apply perhaps only to models with exponential patience. A recent paper of Mandelbaum and Zeltyn [22] addresses this question for (4.1), demonstrating that the linear relation is practically rather robust. See also [23] where we demonstrate a similar linear relation on another data set.

5 Parameter estimation in a call center environment

In order to apply Erlang-A, it is necessary to input values for its four parameters: λ , μ , θ and n . Typical applications use estimates, which are based on historical ACD data, and here we briefly outline procedures for their estimation/prediction. For a more detailed exposition, including some subtleties that occur in practice, readers are referred to [8].

Arrivals: Arrivals of incoming calls are typically assumed Poisson, with time-varying arrival rates. The goal is to estimate/predict these arrival rates, over short time-intervals (15, 30 minutes or one hour), chosen so that the rates are approximately constant during an interval. Then the time-homogeneous model is applied separately over each such interval.

The goal can be achieved in two stages. First, time-series algorithms are used to predict *daily* volumes, taking into account trends and special days (eg. holidays, “Mondays”, special sales). Second, one uses (non)parametric regression techniques for predicting the *fraction* of arrivals per time-interval, out of the daily-total. This fraction, combined with the daily total,

yields actual arrival rates per each time-interval. (See Section 4 in [8] for a detailed treatment.)

Services: Service durations are assumed exponential. Average service times tend to be relatively stable from day to day and from hour to hour. (However, they often change depending on the time-of-day! See [8].)

In practice, service consists of several phases, mainly talk time, wrap-up time (after-call work), and what is sometimes referred to as auxiliary time. An easier-to-grasp notion is thus “idle-time”, namely the time that an agent is immediately accessible for service. It is thus also possible to estimate the average service time during a time interval by:

$$\frac{\text{Total Working Time} - \text{Total Idle Time}}{\text{Number of Served Customers}},$$

where Total Working Time is the product of the Number of Agents by the Interval Duration. (See Adler et al. [1] for an application of this approach, in the context of product development.)

Number of agents: In performance analysis, the number of agents n is an Erlang-A input. In staffing decisions, n is typically an output. In both cases, n is in fact the needed number of FTE’s (Full Time Equivalent positions), and hence it must be normalized by the *rostered staff factor (RSF)*, or *shrinkage factor*, which accounts for absenteeism, unscheduled breaks etc. (See Cleveland and Mayben [10]). For example, if 100 agents are required for answering calls, in fact more agents (105, 110, ...) should be assigned to shift, depending on RSF.

Patience: (Im)patience time is assumed exponential, say $exp(\theta)$. One must then estimate the individual abandonment rate θ , or equivalently, the average patience ($1/\theta$). A difficulty arises from the fact that direct observations are censored - indeed, one can only measure the patience of customers who abandon the system before their service began. For the customers receiving service, their waiting time in queue is only a lower bound for their patience. There are statistical methods for “un-censoring” data; see [8]. Another, more basic problem for estimating θ , is that most ACD data contain only averages, as opposed to call-by-call statistics that are required by the available “uncensoring” methods. To this end, we suggest here two methods for estimating average patience. The first is based on the relation (4.1) between the probability to abandon and average wait. The average wait in queue, $E[W]$, and the fraction of customers abandoning, $P\{\text{Ab}\}$, are in fact standard ACD data outputs, thus, providing the means for estimating θ as follows:

$$\hat{\theta} = \frac{P\{\text{Ab}\}}{E[W]} = \frac{\% \text{Abandonment}}{\text{Average Wait}}.$$

A second more general approach is to calculate some performance measure (see Section 4) and compare the result to the value derived from ACD data. (This approach is applied in [8].) The goal is to calibrate the patience parameter until these estimates closely match. One advantage of this method is the flexibility in choosing the performance measure being matched,

which might depend on the given ACD data. Furthermore, this calibration represents a form of validation of the model’s assumptions, and can compensate for discrepancies.

6 Approximations

Although exact formulae for the Erlang-A system are available and can be incorporated in software (see Sections 3 and 4), they are too complicated for providing guidelines and insights to call center researchers and managers. Consequently, some useful and insightful approximations have been developed, which we now describe.

It has been found useful to distinguish three operational regimes, as in Garnett et al. [15] and Zeltyn and Mandelbaum [23]. Each regime represents a different philosophy for operating a call center. One regime is Efficiency-Driven (ED), another is service Quality-Driven (QD), and the third one rationalizes efficiency and quality, namely it is Quality and Efficiency-Driven (QED).

We are interested mainly in not-too-small call centers. Hence, we think of the service and abandonment rates, μ and θ , as fixed, and the arrival rate λ is large enough (formally, it increases indefinitely).

Actually, the regimes are determined by the *offered load* parameter R , which is defined as $R = \frac{\lambda}{\mu}$; R represents the amount of work, measured in time-units of service, that arrives to the system per unit of time. ($R = \lambda \cdot \frac{1}{\mu}$ is thus a more telling representation.) R is also the staffing level n that would be prescribed by the deterministic approach (see the end of Section 2), which ignores stochastic variability. An emphasis of efficiency (service quality) would conceivably lead to $n < R$ ($n > R$); the deviation of n from R then increases with the intensity of the emphasis. We now proceed with a formal description of the three operational regimes.

QED (Quality and Efficiency-Driven) regime.

$$n \approx R + \beta\sqrt{R}, \quad -\infty < \beta < \infty; \quad (6.1)$$

β is a service-grade parameter – the larger it is, the better is the service-level. The staffing regime, described by (6.1), is governed by the so-called *Square Root Rule*. This rule was already described by Erlang [11], as early as 1924. He reported that it had been in use at the Copenhagen Telephone Company since 1913. A formal QED analysis for the Erlang-C queue appeared only in 1981, in the seminal paper of Halfin and Whitt [16]. (The service grade β must be positive in this case.) Garnett et al. [15] explored Erlang-A in the QED regime and Zeltyn and Mandelbaum [23] treated the M/M/n+G queue with a general patience distribution. (With abandonment, β can be also 0 or negative.)

In the QED regime, the delay probability $P\{W > 0\}$ converges to a constant that is a function of the *service grade* β and the ratio μ/θ (or $\frac{1}{\theta/\mu}$, which is average patience that is measured in units of average service time). The probability to abandon and average wait vanish, as $\lambda, n \uparrow \infty$, at rate $\frac{1}{\sqrt{n}}$. Formulae for different performance measures can be found in [15] or [23]. It is significant and useful to mention that the QED approximations are often valid over a wide range of parameters, from the very large call centers (1000's of agents) to the moderate-size ones (10's of agents).

Figure 7: **Asymptotic relations between service grade and delay probability**

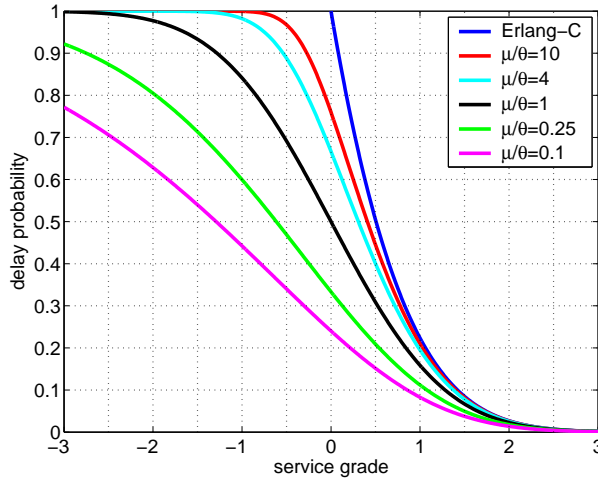


Figure 7 illustrates the dependence between β and $P\{W > 0\}$, for varying values of the ratio μ/θ . In addition, we plotted the curve for the Erlang-C queue, which is meaningful for positive β only. Note that for large values of μ/θ (very patient customers) the Erlang-A curves get close to the Erlang-C curve.

Remark. When $\beta = 0$ in (6.1), the staffing level corresponds to the simple rule that does not take into account stochastic considerations: assign the number of agents equal to the *offered load* λ/μ . In Erlang-C, this “naive” approach would lead to system instability. However, in Erlang-A (which is a much better fit to the real world of call centers than Erlang-C) one would get a reasonable-to-good performance level. For example, if the service rate μ is equal to the individual abandonment rate θ , and $\beta = 0$, 50% of customers would get service immediately upon arrival. (Check it in Figure 7. Note that for Erlang-C, 50% delay probability corresponds to $\beta = 0.5$.) This suggests why some call centers that are managed using simplified deterministic models, actually perform at reasonable service levels. (One obtains the “right answer” from the “wrong reasons”.)

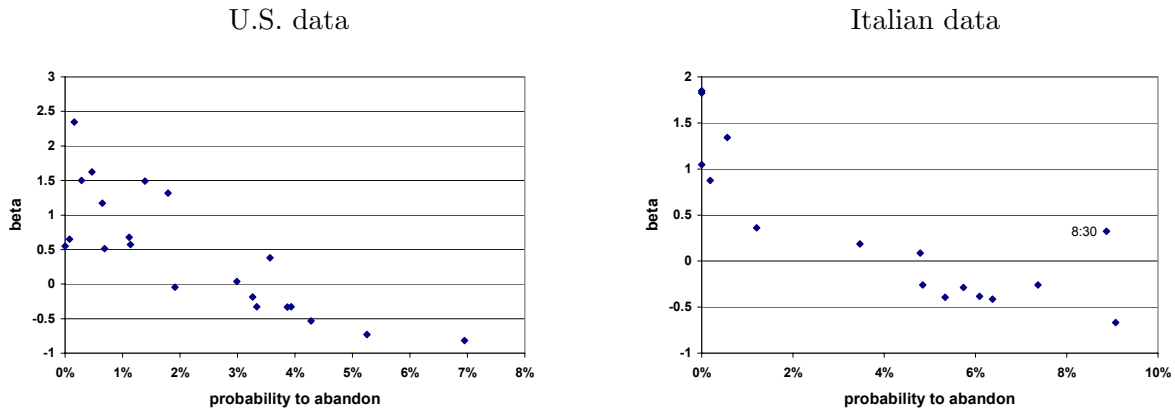
The QED regime enables one to combine high levels of efficiency (agents’ utilization close to 100%) and service quality (agents’ accessibility) . The scatterplots in Figure 8 illustrate this

point. The plots display data from ACD reports of two call centers: Italian and American, collected in half-hour intervals during a single working day. The service grade β is calculated via

$$\beta = \frac{n - R}{\sqrt{R}}.$$

We observe moderate-to-small values of abandonment for the service grade $-1 \leq \beta \leq 2$. Plots of average waiting time exhibit a similar behavior – see [25].

Figure 8: **Service grade for call centers - correlation with abandonment**



ED (Efficiency-Driven) regime.

$$n \approx R \cdot (1 - \gamma), \quad \gamma > 0. \quad (6.2)$$

In this case, virtually all customers wait, the probability to abandon converges to γ and average wait is close to γ/θ . (See Mandelbaum and Zeltyn [24] for additional performance measures.) This regime could be used if efficiency considerations are of main significance. Indeed, it has gained importance in recent research (see, for example, few papers of Whitt [30, 31]), following the observation that ED could yield performance that is acceptable for many call centers, for example those operating in not-for-profit environments.

QD (Quality-Driven) regime.

$$n \approx R \cdot (1 + \gamma), \quad \gamma > 0. \quad (6.3)$$

The staffing regime (6.3) should be implemented if quality considerations far dominate efficiency considerations (e.g. high-valued customers or emergency phones). Major performance measures (delay probability, fraction abandoning, average wait) vanish here at an exponential rate of n .

Remark. Above, we considered steady-state performance measures. Process-limit results for the number-in-system process $L = \{L(t), t \geq 0\}$ are available for the QED and ED regimes in Garnett et al. [15] and Whitt [30], respectively.

7 Applications to call centers

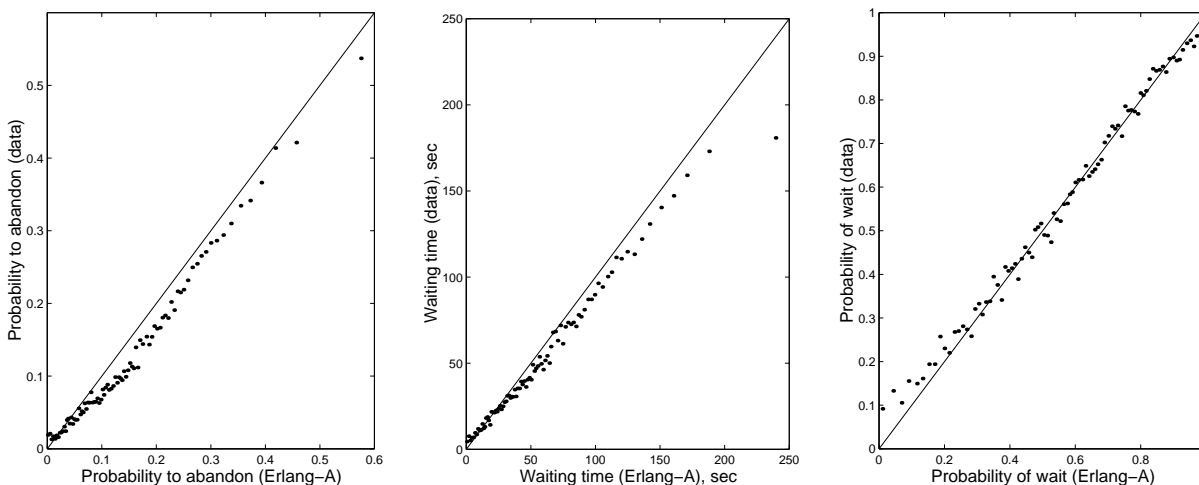
7.1 Erlang-A performance measures: comparison against real data

We now validate the Erlang-A model against the hourly data for the Israeli bank call center, already used for the example in Section 4. Three performance measures are considered: probability to abandon, average waiting time and probability of wait. Their values are calculated for the hourly intervals using exact Erlang-A formulae. Then the results are aggregated along the same method employed in Figure 6. The resulting 86 points are compared against the line $y = x$: the better the fit the better Erlang-A describes reality.

Computation of the Erlang-A parameters. Parameters λ and μ are calculated for every hourly interval. We also calculate each hour's average number of agents n . Because the resulting n 's need not be integral, we apply a continuous extrapolation of the Erlang-A formulae, obtained from relationships developed in [26]. Finally, for θ we use formula (4.1).

The results are displayed in Figure 9. The figure's two left-hand graphs exhibit a relatively small yet consistent overestimation with respect to empirical values, for moderately and highly loaded hours. The right-hand graph shows a very good fit everywhere, except for very lightly and very heavily loaded hours. The underestimation for small values of $P\{W > 0\}$ can be probably attributed to violations of work conservation (idle agents do not always answer a call immediately). Summarizing, it seems that these Erlang-A estimates can be used as close *upper bounds* for the main performance characteristics of our call center.

Figure 9: Erlang-A formulas vs. data averages [8]



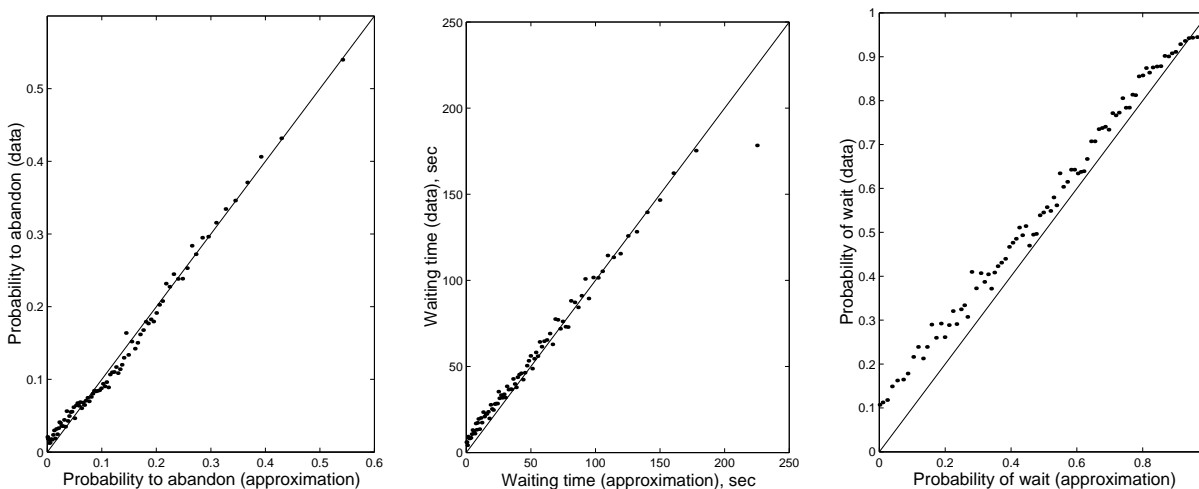
7.2 Erlang-A approximations: comparison against real data

In Section 6 we discussed approximations of various performance measures for the Erlang-A (M/M/n+M) model. Such approximations require significantly less computational effort than exact Erlang-A formulae. Figure 10, based on the same data as Figure 9, demonstrates a good fit between data averages and the approximations.

In fact, the fits for the probability of abandonment and average waiting time are somewhat superior to those in Figure 9 (the approximations provide somewhat larger values than the exact formulae). This phenomenon suggests two interrelated research questions of interest: explaining the overestimation in Figure 9 and better understanding the relationship between Erlang-A formulae and their approximations.

The empirical fit of the simple Erlang-A model and its approximation turns out to be very (perhaps surprisingly) accurate. Thus, for the call center in consideration – and those like it – use of Erlang-A for workforce management could and should improve operational performance.

Figure 10: Erlang-A approximations vs. data averages [8]



8 Some advanced features of 4CallCenters

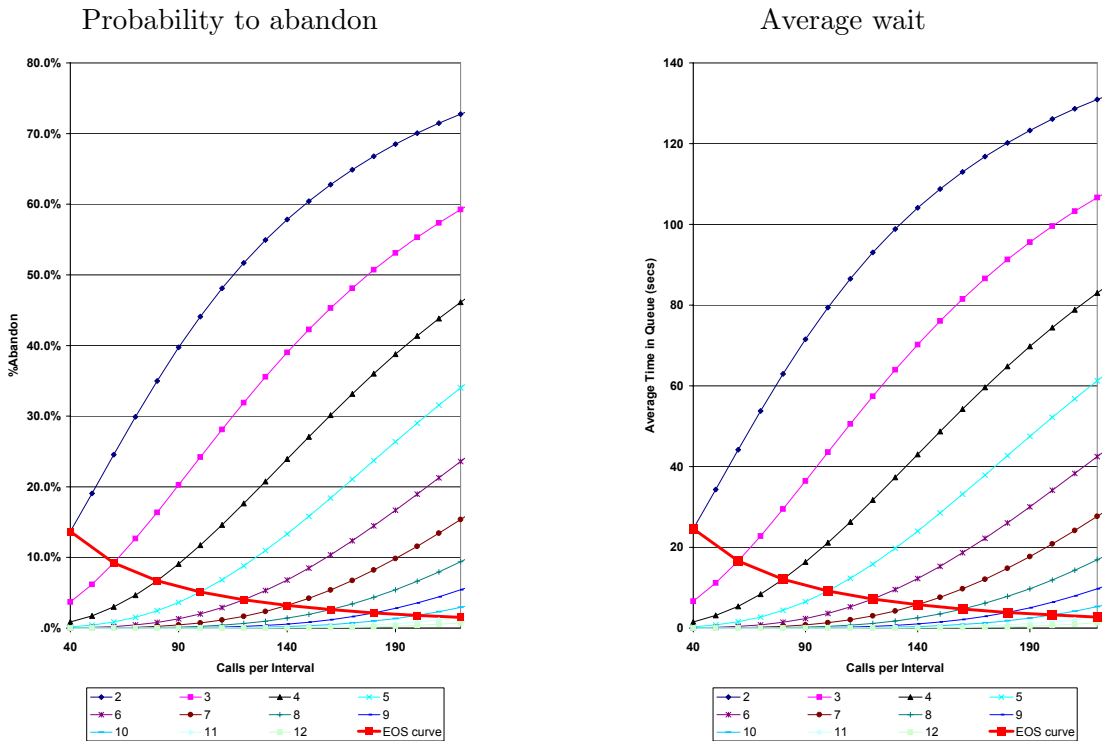
The 4CallCenters software [12] provides a valuable tool for implementing Erlang-A calculations. Its basic feature is “Performance Profiler” that enables calculation of all the useful performance measures, given the four Erlang-A parameters as input. In addition, 4CallCenters allows many advanced options: staffing queries, graphs, export and import of data and more.

Here we demonstrate, as an example, two advanced capabilities of 4CallCenters.

Example 1: Advanced profiling. One can vary any input parameters of the Erlang-A queue and display the corresponding model output (performance measures) either in a table or graphically. For example, let the average service time equal 2 minutes and average patience 3 minutes. Let the arrival rate vary from 40 to 230 calls per hour, in steps of 10, and the number of agents from 2 to 12. Then one can immediately produce a table that contains values of different performance measures for all combinations of the two input parameters.

Figure 11 shows the dependence of the probability to abandon and average wait on different number of agents. Note that the two plots look identical: the reason is relation (4.1). In addition, the red curves on both plots in Figure 11 illustrate Economies of Scale (EOS): while offered load per server remains constant along this curve $\left(\frac{\lambda}{n\mu} = \frac{2}{3}\right)$, performance significantly improves as the number of agents increases. For example, the probability to abandon is equal to 13.7% for $n = 2$, 5.1% for $n = 5$ and 1.5% for $n = 12$. Finally, note that both $P\{Ab\}$ and $E[W]$ actually vanish as n gets large.

Figure 11: 4CallCenters. Advanced profiling.



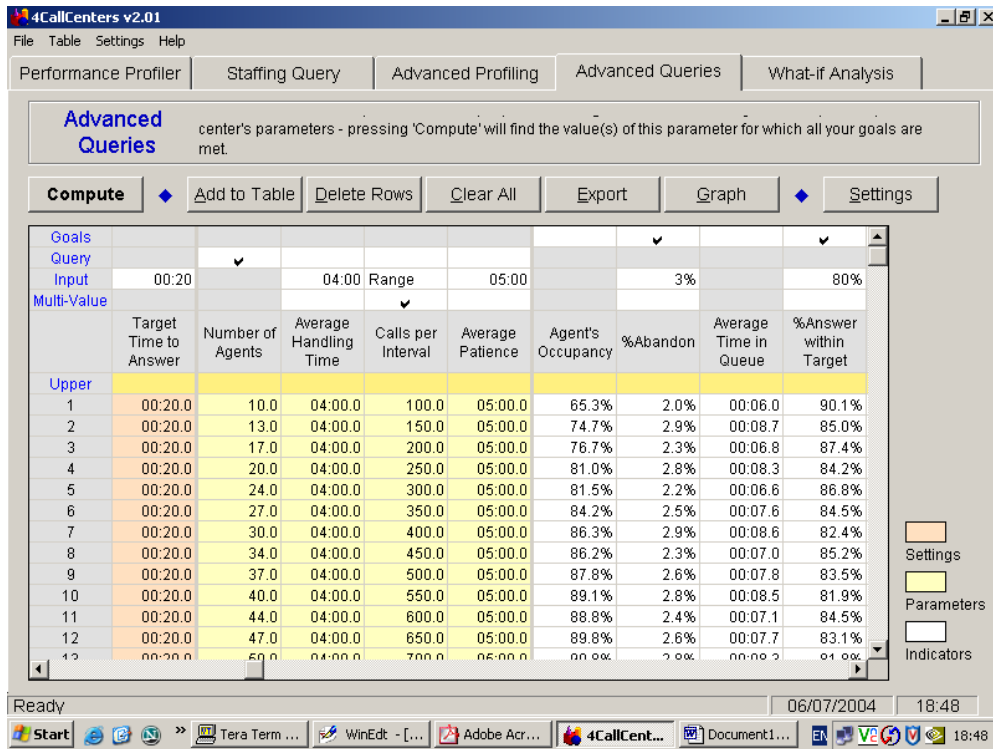
Example 2: Advanced staffing queries. 4CallCenters enables staffing queries with several performance goals. For example, assume that the average service time is equal to 4 minutes, and average patience is 5 minutes. Our goal is to calculate appropriate staffing levels for arrival-rate values that vary from 100 to 1200, in steps of 50. The performance targets are:

- Probability to abandon less than 3%;
- 80% of customers served within 20 seconds.

Figure 12 presents the screen output of 4CallCenters.

The first plot of Figure 13 displays the minimal staffing level that adheres to both goals. The EOS phenomenon is observed here as well: 10 agents are needed for 100 calls per hour but only 83 (rather than $10 \cdot 12 = 120$) for 1200 calls per hour. (Despite its look, the curve in the first plot is not a straight line.) The second plot displays the values of the two target performance measures. (This plot, unlike the first one, is not an immediate output graph of 4CallCenters but rather an edited version of it.)

Figure 12: 4Callcenters. Advanced staffing queries.



Remark. Since the number of agents must be an integer, we observe performance “zigzags” in the right plot of Figure 13.

9 Some open research topics

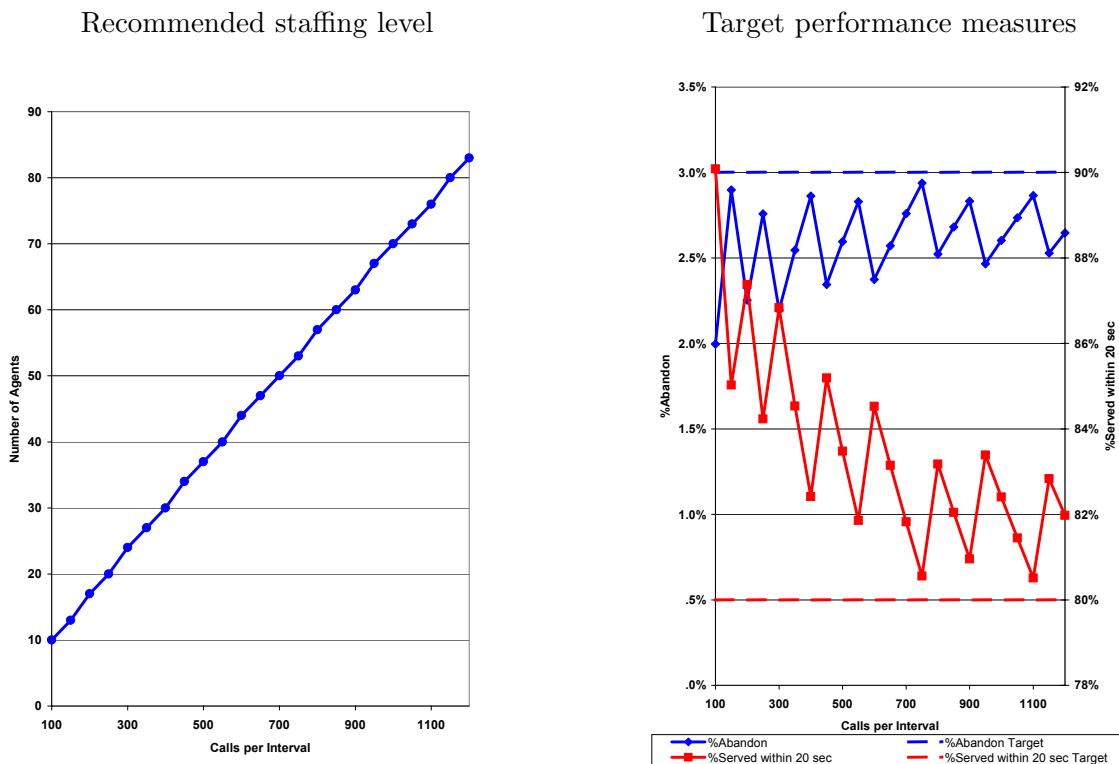
9.1 Dimensioning the Erlang-A queue

One can search for an optimal staffing level, given the trade-off between staffing cost, cost of customers' waiting and cost of abandonment. Borst et al. [6] referred to this problem as *dimensioning* and solved it for the Erlang-C queue (no abandonment): if the staffing cost and the cost of waiting are comparable, the optimal staffing should take place in the QED regime, described in Section 6 (with positive β in formula (6.1)). Ongoing research by Borst et al. [7] is dedicated to the same question for the Erlang-A queue. For example, let the average operational cost (per unit of time) be equal to

$$U(n, \lambda) = c \cdot n + \lambda a \cdot P\{\text{Ab}\},$$

where c is the staffing cost, and a is the abandonment cost. Our goal is to minimize cost. (Note that this is in fact mathematically equivalent to maximizing revenues.)

Figure 13: 4CallCenters. Dynamics of staffing level and performance.



Define the abandonment/staffing cost ratio by $r \triangleq a/c$, and let $s \triangleq \sqrt{\mu/\theta}$. Assume that $a > c/\mu$. (Otherwise, the asymptotic optimal policy is $n^* = 0$: not to provide service at all.) Then we suggest that the asymptotic optimal staffing level is equal to

$$n^* = \lceil R + y^*(r; s) \cdot \sqrt{R} \rceil, \quad (9.1)$$

where the square brackets in (9.1) denote the nearest integer value and the function $y^*(\cdot)$ is defined by

$$y^*(r; s) \triangleq \arg \min_{-\infty < y < \infty} \left\{ c \cdot y + a \cdot \theta s \cdot \left[1 + \frac{h(ys)}{sh(-y)} \right]^{-1} \cdot [h(ys) - ys] \right\},$$

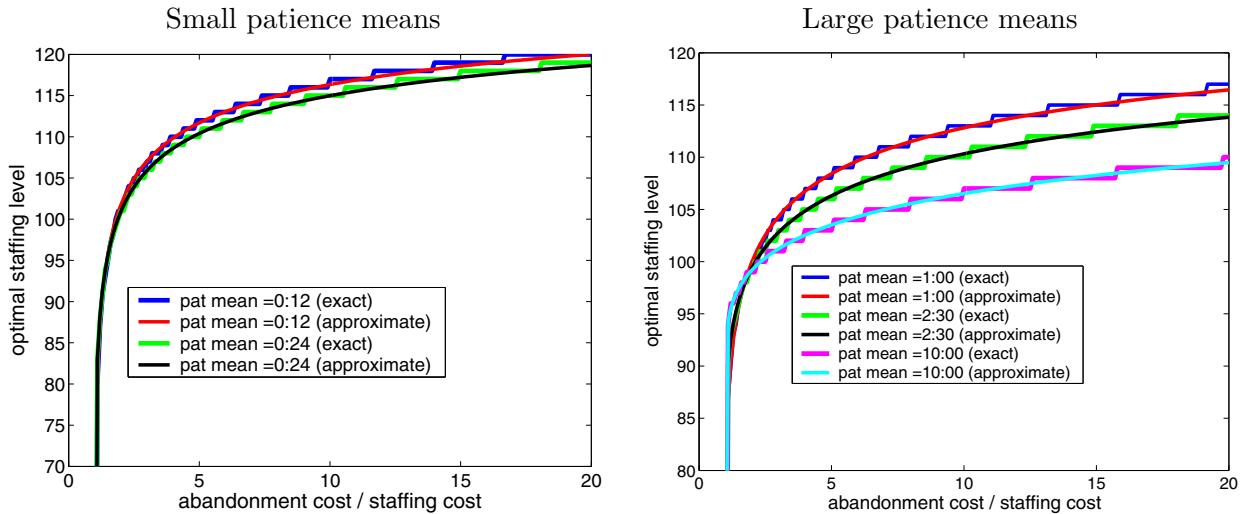
and $h(\cdot) = \phi(\cdot)/(1 - \Phi(\cdot))$ is the hazard rate of the standard normal distribution ($\phi(\cdot)$ is its density function and $\Phi(\cdot)$ is the cumulative distribution function).

As in [6], Figure 14 compares the rule (9.1) with the exact optimal staffing values. We consider five exponential patience distributions with different means and perform comparisons by varying the value of the ratio r . A perfect fit is observed for *all* the special cases!

Numerical experiments for other cost optimization problems (e.g. with waiting cost, instead of abandonment cost) also demonstrate a very close correspondence between exact values and the corresponding analogs of (9.1). Hence the goal is to develop a theoretical framework, parallel to [6], that will support our experimental research. In addition, we are working on a *constraint satisfaction* version where one chooses the least number of agents that adheres to a given constraint on the waiting and/or the abandonment cost. This latter formulation is in fact closer to the way that managers perceive their staffing problems in practice.

Figure 14: **Cost optimization. Approximation vs. exact optimum.**

Arrival rate $\lambda = 100$, service rate $\mu = 1$.



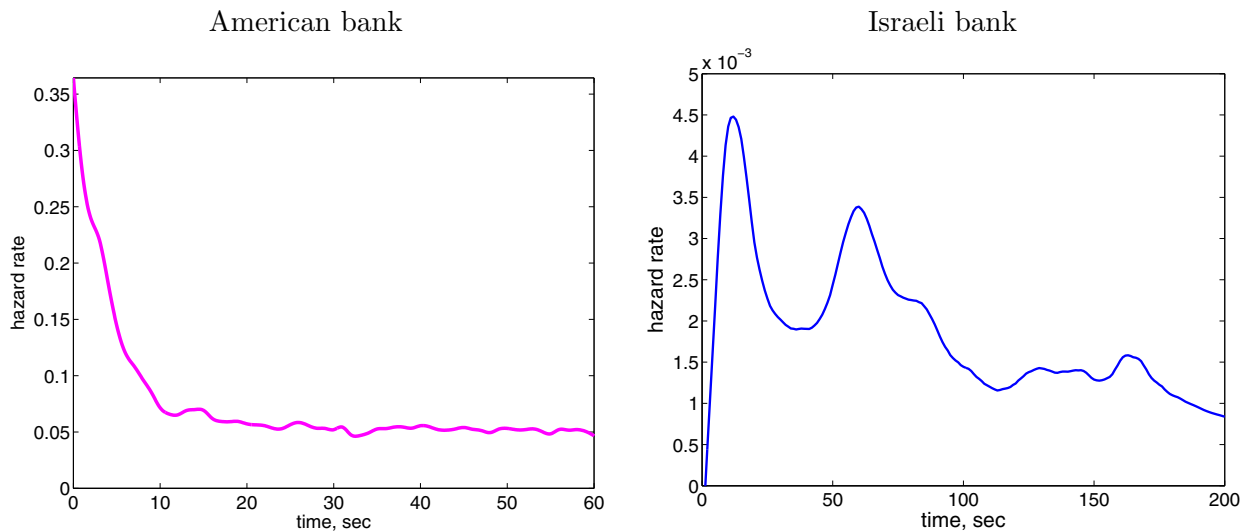
9.2 Human behavior

The Erlang-A model assumes exponential iid patience times that do not depend on the state of the system, time-of-day etc. In practice, these assumptions are not always valid.

In Figure 15 we display estimates of the hazard rates of the customers' patience for two banks: a large U.S. bank and a small Israeli one. In the two cases we observe different, but clearly non-exponential patterns. (Recall that the hazard rate of an exponential random variable is a constant.) American customers are very impatient at the beginning of their wait, but their patience stabilizes after approximately 10 seconds. In contrast, Israeli customers have two clear peaks of abandonment: approximately at 15 and at 60 seconds. (It turns out that these two surges of abandonment take place immediately after two recorded messages to which customers are exposed: the first one when they enter the queue and the second after approximately 1 minute.)

Therefore, at least in some applications customers' patience times are non-exponential and applicability of the Erlang-A formulae to such systems should be studied. (Recall Section 7.)

Figure 15: **Bank data: hazard rates of patience times**



Patience index. In search for a better understanding of customers' (im)patience, we have found a relative definition to be of use. Specifically, we define the *patience index* to be

$$\begin{aligned} \text{Theoretical Patience Index} &\triangleq \frac{\text{time a customer is willing to wait}}{\text{time a customer is required to wait}} \\ &= \frac{\text{average patience}}{\text{average offered wait}} = \frac{E[\tau]}{E[V]}. \end{aligned}$$

While this patience index makes sense intuitively, its calculation requires the application of survival analysis techniques to call-by-call data. Such data may not be available in certain circumstances. Therefore, we wish to find an *empirical* index which will work as an auxiliary measure for the patience index.

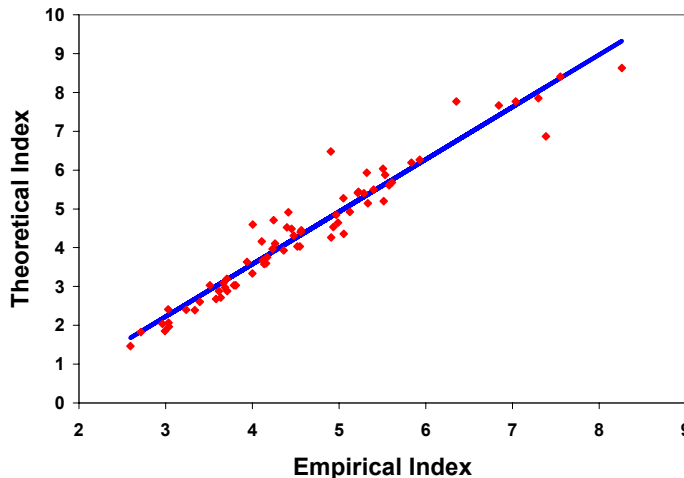
We found the following to be a very useful definition:

$$\text{Empirical Patience Index} \triangleq \frac{\% \text{ served}}{\% \text{ abandoned}}.$$

The empirical index is easily calculable since both the numbers of served and of abandoned calls are very easy to obtain from call-center reports.

Figure 16 demonstrates how well the empirical patience index estimates the theoretical patience index for the Israeli bank data [8]. (Aggregated data of 68 quarter hours between 7am and midnight is used.)

Figure 16: **Patience index – empirical vs. theoretical** [9]



Under certain circumstances, one can explain the closeness of the theoretical and empirical indices [8]. However, these explanations are unsatisfactory and hence leave open this research direction.

Adaptive behavior. In the papers [21, 29, 35], Mandelbaum, Shimkin and Zohar analyze models for an *adaptive* behavior of tele-customers, which “tune” their patience according to anticipated or perceived systems congestion. Data from the call center of our Israeli bank supports the applicability of these models, but more is to be done in this direction.

9.3 Uncertainty in parameter values

In the real world, one never knows the exact values of the four Erlang-A parameters. Therefore, it is essential to study the sensitivity of performance measures. In a recent paper [33], Whitt calculates *elasticities* in Erlang-A, which measure the percentage change of a performance measure caused by a small percentage change in a parameter. Both exact numerical algorithm and several types of approximations are used. It turns out that Erlang-A performance is quite sensitive to small changes in the arrival rate, service rate, or number of agents, but relatively insensitive to small changes in the abandonment rate.

In staffing planning, of which the number of agents is the output, it is reasonable to assume knowledge of the service rate. Thus, the problem of uncertainty in the arrival rate surfaces as the most significant.

In Brown et al. [8] it was shown that the Poisson arrival rate in an Israeli call center varies from day to day and its prediction raised statistical and practical challenges. This motivates the study of queueing models, in which the Poisson arrival rate Λ (the arrival-rate function) is a random variable (random process).

If $E(\Lambda) \rightarrow \infty$ and its standard deviation is of the order $\sqrt{E(\Lambda)}$, we expect that the QED operational regime and the square-root staffing rule will play a role that is similar to the one with known (deterministic) arrival rate; the offered load R in (6.1) will be replaced by the average offered load $E(\Lambda)/\mu$, and uncertainty will manifest itself through a different value of the service grade β . However, if $\sigma(\Lambda)$ is of the order $E(\Lambda)$, the “cruder” ED regime seems to be the most appropriate; see Whitt [32], and Bassamboo, Harrison and Zeevi [4].

References

- [1] 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. 5
- [2] Aguir M.S., Karaesmen F., Aksin O.Z. and Chauvet F. (2004) The impact of retrials on call center performance, *OR Spectrum*, Special Issue on Call Center Management, 26(3),353-376. 1
- [3] Bain P. and Taylor P. (2002) Consolidation, “Cowboys” and the developing employment relationship in British, Dutch and US call centres. In: Holtgrewe U., Kerst C. and Shire K. (Ed.), *Re-Organising Service Work*. Ashgate Publishing Limited, 42-62. 1
- [4] Bassamboo A., Harrison J.M. and Zeevi A. (2004) Design and Control of a Large Call Center: Asymptotic Analysis of an LP-based Method. Submitted for publication. 9.3

- [5] Bittner S., Schietinger M., Schroth J. and Weinkopf C. (2002) Call Centres in Germany: Employment, Training and Job Design. In: Holtgrewe U., Kerst C. and Shire K. (Ed.), *Re-Organising Service Work*. Ashgate Publishing Limited, 63-85. [1](#)
- [6] Borst S., Mandelbaum A., and Reiman M. (2004), Dimensioning large call centers, *Operations Research*, 52(1), 17-34. [9.1](#), [9.1](#)
- [7] Borst S., Mandelbaum A., Reiman M. and Zeltyn S. (2004) Dimensioning call centers with abandonment. In prepatation. [9.1](#)
- [8] Brown L.D., Gans N., Mandelbaum A., Sakov A., Shen H., Zeltyn S. and Zhao L. (2002) Statistical analysis of a telephone call center: a queueing science perspective. To be published in *JASA*. [1](#), [6](#), [4.4](#), [5](#), [9](#), [10](#), [9.2](#), [9.2](#), [9.3](#)
- [9] Call Center Data (2002) Technion, Israel Institute of Technology. Available at <http://iew3.technion.ac.il/serveng/callcenterdata/index.html>. [4.4](#), [16](#)
- [10] Cleveland B., Mayben J. (1997) *Call Center Management on Fast Forward*. Annapolis: Call Center Press. [5](#)
- [11] Erlang A.K. (1948) On the rational determination of the number of circuits. In *The life and works of A.K.Erlang*. Brockmeyer E., Halstrom H.L. and Jensen A., eds. Copenhagen: The Copenhagen Telephone Company. [6](#)
- [12] 4CallCenters Software (2002). Available at <http://iew3.technion.ac.il/serveng/4CallCenters/Downloads.htm>. [4.1](#), [4.2](#), [8](#), [A](#)
- [13] Gans N., Koole G. and Mandelbaum A. (2003) Telephone call centers: a tutorial and literature review. Invited review paper, *Manufacturing and Service Operations Management*, 5(2), 79-141. [1](#)
- [14] Garnett O. and Mandelbaum A. (2000) An Introduction to Skills-Based Routing and its Operational Complexities. Teaching note, Technion, Israel. Available at <http://iew3.technion.ac.il/serveng2004/Lectures/SBR.pdf>. [2](#)
- [15] Garnett O., Mandelbaum A. and Reiman M. (2002) Designing a telephone call-center with impatient customers. *Manufacturing and Service Operations Management*, 4,208-227. [1](#), [2](#), [4.2](#), [6](#), [6](#), [6](#)
- [16] Halfin S. and Whitt W. (1981) Heavy-traffic limits for queues with many exponential servers. *Operations Research*, 29, 567-588. [6](#)
- [17] Helber S. and Mandelbaum A. (2004) GIF Research Proposal. [1](#)

- [18] Helber S. and Stolletz R. (2004) *Call Center Management in der Praxis*. Springer-Verlag, Berlin, Heidelberg. (In German) [1](#)
- [19] Mandelbaum A., Sakov A. and Zeltyn S. (2001) Empirical analysis of a call center. Technical report, Technion. Available at <http://iew3.technion.ac.il/serveng/References/references.html>. [4.4](#)
- [20] Mandelbaum A. (2003) Call Centers. Research Bibliography with Abstracts. Version 5. Available at <http://iew3.technion.ac.il/serveng/References/references.html>. [1](#)
- [21] Mandelbaum A. and Shimkin N. (2000) A model for rational abandonment from invisible queues. *Queueing Systems: Theory and Applications (QUESTA)*, 36, 141-173. [9.2](#)
- [22] Mandelbaum A. and Zeltyn S. (2004) The Impact of Customers Patience on Delay and Abandonment: Some Empirically-Driven Experiments with the M/M/N+G Queue. *OR Spectrum*, Special Issue on Call Center Management, 26(3), 377-411. [1](#), [4.4](#)
- [23] Mandelbaum A. and Zeltyn S. (2004) Call centers with impatient customers: many-server asymptotics of the M/M/n+G queue. Submitted to *QUESTA*. Available at <http://iew3.technion.ac.il/serveng/References/references.html>. [4.4](#), [6](#), [6](#)
- [24] Mandelbaum A. and Zeltyn S. (2004) M/M/n+G queue. Summary of performance measures. Available at <http://iew3.technion.ac.il/serveng/References/references.html>. [4.3](#), [6](#), [A](#)
- [25] Mandelbaum A. and Zeltyn S. (2004) The Palm/Erlang-A Queue, with Applications to Call Centers. Teaching note to *Service Engineering* course. Available at <http://iew3.technion.ac.il/serveng/References/references.html>. [6](#)
- [26] Palm C. (1957) Research on telephone traffic carried by full availability groups. *Tele*, vol.1, 107 pp. (English translation of results first published in 1946 in Swedish in the same journal, which was then entitled *Tekniska Meddelanden fran Kungl. Telegrafstyrelsen*.) [2](#), [7.1](#), [A](#), [A](#)
- [27] Riordan J. (1962) *Stochastic Service Systems*, Wiley. [4.3](#)
- [28] “Service Engineering” course web-site, Technion, <http://iew3.technion.ac.il/serveng>. [4.4](#)
- [29] Shimkin N. and Mandelbaum A. (2004) Rational abandonment from tele-queues: non-linear waiting costs with heterogeneous preferences. *Queueing Systems: Theory and Applications (QUESTA)*, 47, 117-146. [9.2](#)

- [30] Whitt W. (2004) Fluid Models for Many-Server Queues with Abandonments. Submitted to *Operations Research*. **6, 6**
- [31] Whitt W. (2004) Two Fluid Approximations for Multi-Server Queues with Abandonments. Submitted to *Operations Research Letters*. **6**
- [32] Whitt W. (2004) Staffing a Call Center with Uncertain Arrival Rate and Absenteeism. Submitted to *Management Science*. **9.3**
- [33] Whitt W. (2004) Sensitivity of Performance in the Erlang A Model to Changes in the Model Parameters. Submitted to *Operations Research*. **9.3**
- [34] Zeltyn S. (2004) Call centers with impatient customers: exact analysis and many-server asymptotics of the M/M/n+G queue, Ph.D. Thesis, Technion. Available at <http://iew3.technion.ac.il/serveng/References/references.html>. **2**
- [35] Zohar E., Mandelbaum A. and Shimkin N. (2002) Adaptive behavior of impatient customers in tele-queues: theory and empirical support. *Management Science*, 48, 566-583. **9.2**

A The Erlang-A queue: useful formulae for the steady-state distribution and some performance measures

Steady-state distribution. Palm [26] derived the following representation for the steady-state distribution, defined in (3.1):

$$\pi_j = \begin{cases} \pi_n \cdot \frac{n!}{j! \cdot (\lambda/\mu)^{n-j}}, & 0 \leq j \leq n, \\ \pi_n \cdot \frac{(\lambda/\theta)^{j-n}}{\prod_{k=1}^{j-n} (\frac{n\mu}{\theta} + k)}, & j \geq n+1, \end{cases}$$

where

$$\pi_n = \frac{E_{1,n}}{1 + \left[A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) - 1 \right] \cdot E_{1,n}},$$

$$A(x, y) \triangleq \frac{x e^y}{y^x} \cdot \gamma(x, y),$$

and

$$\gamma(x, y) \triangleq \int_0^y t^{x-1} e^{-t} dt, \quad x > 0, y \geq 0.$$

is the *incomplete Gamma function*; $E_{1,n}$ denotes the *blocking probability* in the M/M/n/n (Erlang-B) system:

$$E_{1,n} = \frac{\frac{(\lambda/\mu)^n}{n!}}{\sum_{j=0}^n \frac{(\lambda/\mu)^j}{j!}}.$$

Remark. A simple way for calculating $E_{1,n}$ is the recursion

$$E_{1,0} = 0; \quad E_{1,n} = \frac{\rho E_{1,n-1}}{1 + \rho E_{1,n-1}}, \quad n \geq 1,$$

in which ρ is the offered load per agent, namely

$$\rho \triangleq \frac{\lambda}{n\mu}.$$

Performance measures. As discussed above, 4CallCenters [12] provides a convenient tool for Erlang-A calculations. Mandelbaum and Zeltyn [24] present a theoretical framework for computations in the more general M/M/n+G system and explain how to adapt it to Erlang-A. Both approaches can be used for calculation of the practical measures from Subsection 4.1. Below we present several explicit expressions. Average wait and the probability to abandon are the most widely used performance measures in practice. To these we add the probability of wait, which is important in view of the fact that it characterizes the operational regime (ED, QD or QED) – recall Section 6.

Probability of wait. Following Palm [26],

$$P\{W > 0\} = \frac{A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) \cdot E_{1,n}}{1 + \left(A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) - 1\right) \cdot E_{1,n}}, \quad (\text{A.1})$$

Probability to abandon. The probability to abandon of delayed customers is equal to

$$P[\text{Ab}|W > 0] = \frac{1}{\rho A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right)} + 1 - \frac{1}{\rho}. \quad (\text{A.2})$$

The fraction abandoning, $P\{\text{Ab}\}$, is simply the product $P[\text{Ab}|W > 0] \times P\{W > 0\}$.

Average waiting time. Average waiting time of delayed customers is computed via (A.2) and (4.1):

$$E[W|W > 0] = \frac{1}{\theta} \cdot \left[\frac{1}{\rho A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right)} + 1 - \frac{1}{\rho} \right]. \quad (\text{A.3})$$

The unconditional average wait $E[W]$ equals the product of (A.1) with (A.3).