

Dimensioning Large Call Centers

Sem Borst

CWI, P. O. Box 94079, 1090 GB Amsterdam, The Netherlands, and Bell Labs, Lucent Technologies, Murray Hill, New Jersey 07974-0636, sem.borst@cwi.nl

Avi Mandelbaum

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

Martin I. Reiman

Bell Labs, Lucent Technologies, Murray Hill, New Jersey 07974-0636, marty@research.bell-labs.com

We develop a framework for asymptotic optimization of a queueing system. The motivation is the staffing problem of large call centers, which we have modeled as M/M/N queues with N , the number of agents, being large. Within our framework, we determine the asymptotically optimal staffing level N^* that trades off agents' costs with service quality: the higher the latter, the more expensive is the former. As an alternative to this optimization, we also develop a constraint satisfaction approach where one chooses the least N^* that adheres to a given constraint on waiting cost. Either way, the analysis gives rise to three regimes of operation: quality-driven, where the focus is on service quality; efficiency-driven, which emphasizes agents' costs; and a rationalized regime that balances, and in fact unifies, the other two. Numerical experiments reveal remarkable accuracy of our asymptotic approximations: over a wide range of parameters, from the very small to the extremely large, N^* is *exactly* optimal, or it is accurate to within a single agent. We demonstrate the utility of our approach by revisiting the square-root safety staffing principle, which is a long-existing rule of thumb for staffing the M/M/N queue. In its simplest form, our rule is as follows: if c is the hourly cost of an agent, and a is the hourly cost of customers' delay, then $N^* = R + y^*(a/c)\sqrt{R}$, where R is the offered load, and $y^*(\cdot)$ is a function that is easily computable.

Subject classifications: Queues, optimization: choosing optimal number of servers. Queues, limit theorems: many server queues.

Area of review: Stochastic Models.

History: Received December 2000; revision received May 2002; accepted March 2003.

1. Introduction

Worldwide, telephone-based services have been expanding dramatically in both volume and scope. This has given rise to a huge growth industry—the (telephone) call center industry. Indeed, some assess (Call Center Statistics 2000) that 70% of all customer-business interactions in the U.S. occur in call centers, which employ about 3% of the U.S. workforce (several million agents). Marketing managers refer to call centers as the modern business frontier, being the focus of Customer Relationship Management (CRM); Operations managers are challenged with the fact that personnel costs, specifically staffing, account for over 65% of the cost of running the typical call center. The trade-off between service quality (marketing) and efficiency (operations), thus, naturally arises, and a central goal of ours is to contribute to its understanding.

We argue that call centers typify an emerging business environment in which the traditional quality-efficiency trade-off paradigm could collapse: Extremely high levels of both service quality and efficiency can coexist. Consider, for example, a best-practice U.S. sales call center that attends to an average of 15,000 phone callers daily; the average duration of a call is four minutes, and the variability of calls is significant; agents are highly utilized

(over 90%), yet customers essentially never encounter a busy signal, hardly anyone abandons while waiting, the average wait for service is a mere few seconds, and about *half* of the customers find, upon calling, an idle agent to serve them immediately. Prerequisites for sustaining such performance, to the best of our judgment, are technology-enabled economies of scale and scientifically-based managerial principles and laws. In this paper, we develop an analytical framework (§§4 and 9) that supports such principles. It is based on asymptotic optimization, which yields insight that does not come out of exact analysis. A convincing example is the *square-root safety staffing principle*, described in §2 below. It supports simple, useful rules of thumb for staffing large call centers, rules that so far have been justified only heuristically. Indeed, rigorous asymptotic justifications of such rules are not common in the operations research literature. Hence, another goal here is to convince the reader of their benefits.

1.1. Costs, Optimization, and Constraint Satisfaction

The cost of staffing is the principal component in the operating expenses of a call center. The staffing level is also

the dominant factor to determine service level, as measured in terms of delay statistics: Poor service levels incur either opportunity losses due to deteriorating goodwill, or more direct revenue losses in case of abandonment and blocking (busy signals). While the need to balance service quality and staffing cost is universal, the weight placed on each may vary dramatically. In some call centers, providing maximal customer care is the primary drive whereas in others, handling a high traffic volume at minimal cost is the overriding goal. The challenge, so we argue, is to translate such strategically articulated goals into concrete staffing levels: Simply put, *how many agents are to be staffed in order to provide acceptable service quality and operational efficiency?* In this paper, we answer this question for the M/M/N (Erlang-C) queue, which is the simplest yet most prevalent model that supports call center staffing. In future research we hope to add crucial features of call centers such as abandonment and retrials (Mandelbaum et al. 2002).

Within the M/M/N model, we postulate a staffing cost function $F(N)$ for employing N agents. We assume that (a continuous extension of) $F(N)$ is convex and strictly increasing, which also covers linear costs. The convexity assumption is motivated by the property that the hourly salary tends to increase with the demand in tight labor markets. The fact that the supply of labor is an issue is indirectly supported by the observation that the availability of a low-cost labor force is a major consideration for the location of call center businesses. Low costs (small N) give rise to long waits, which we quantify in terms of a delay cost function $D(t)$ for a customer being served after waiting t units of time. When F dominates D (or conversely D dominates F), the least costs are achieved in an *efficiency-driven* (or conversely a *quality-driven*) operation. When F and D are comparable, optimization leads to a *rationalized* operation which, as it turns out, is robust enough to encompass most circumstances. Formally, the three regimes emerge from an asymptotic analysis of the M/M/N queue, as the arrival rate λ and, accordingly, the optimal staffing level N_λ^* , both scale up to infinity. We refer to such responsive staffing, in response to increased load, as *dimensioning* the call center, which inspired our title. While the staffing levels that we recommend are only asymptotically optimal, they are nevertheless remarkably accurate—to within a *single* agent in the majority of cases. The asymptotics also provide insight, beyond that of exact analysis, about the dependence of the optimal N_λ^* on λ , F , and D .

In industry practice, staffing levels are rarely determined through optimization. One reason is that there is no standard practice for quantifying waiting costs, let alone abandonment, busy-signal and retrial costs; see Andrews and Parsons (1993) for some attempts. Thus, if not by mere experience-based guessing, common practice seeks the least number of agents N^* that satisfies a given constraint on service level. The latter is expressed in terms of some congestion measure, for example the industry-standard Total Service Factor (TSF) given by

$$\text{TSF} = \Pr\{\text{Wait} > T\}, \quad \text{for some } T \geq 0,$$

perhaps combined with 1-800 operating costs. We call this practice *constraint satisfaction*. It is to be contrasted with our previous *optimization* practice, where N^* was determined by cost minimization.

1.2. Introduction to Our Asymptotic Framework

As already mentioned, the justification for our proposed optimal staffing level is based on an asymptotic framework, which we formally develop in §4. Its basic idea is as follows. The primitives of our call center model are the arrival rate λ , the number of servers N , and the average service time $1/\mu$. The latter will be fixed throughout our analysis, while $\lambda \uparrow \infty$ is our asymptotic regime, and N is the parameter over which we optimize. Specifically, given the staffing cost function $F(N)$ and the customer's cost of delay $D(t)$, we express the overall cost per unit of time $C(N, \lambda)$ in terms of three entities: staffing costs, waiting costs, and the probability that an arriving customer is delayed in queue (Erlang-C formula); see (7). Our goal is to solve the discrete optimization problem that seeks N_λ^* which minimizes $C(N, \lambda)$ and, no less importantly, understand the behavior of N_λ^* for large λ . To this end, we translate the discrete optimization problem into a continuous one that is easier to solve, which is carried out by replacing the three entities above with continuous approximations. The optimal solution for the continuous optimization problem provides an approximately optimal solution to our original discrete problem. The approximation is asymptotically optimal in that, as λ increases indefinitely, the ratio of the overall cost at the approximate staffing level to the cost at the true optimal level (both reduced by the cost of staffing at the least level λ/μ needed to assure stability) converges to unity; see Corollary 4.3, and the discussion following it.

Having set up the framework for asymptotic optimality, we then identify continuous approximations to the three cost entities, doing it separately for each of the rationalized, efficiency- and quality-driven regimes described in the previous subsection (see §§6–8, respectively). We then derive similar approximations in the context of constraint satisfaction (§9).

Of central importance to our approximation is the asymptotic analysis of Halfin and Whitt (1981), especially their approximation to the Erlang-C delay function (Lemma 5.1). It gives rise to a square-root safety staffing principle, which reads roughly as follows: For Erlang-C, staffing levels must always exceed the offered load (λ/μ) to ensure stability; this excess is naturally measured in units of the square-root of the offered load, and our optimization problems, in fact, search for the optimal number y^* of such units. The value of y^* depends on the operational regime under discussion. For example, y^* in the rationalized regime is a function of the cost data, which is independent of λ . (The special case with linear staffing and waiting costs is presented in the next section.) On the other hand, in the efficiency-driven regime where fewer resources suffice, y^* vanishes as $\lambda \uparrow \infty$. The fundamental law behind the square-root scaling is the central limit theorem—see Whitt (1992) for further insight.

1.3. Structure of the Paper

The next section is devoted to an exposition of the square-root safety staffing principle, followed by a review of the related literature.

In §3, we set up our M/M/N model and its cost structure. The framework for asymptotic optimality is introduced in §4. Its applications require some special functions which are introduced in §5, notably the Halfin-Whitt delay function $P(\cdot)$ (Halfin and Whitt 1981), plotted in Figure 3. It provides an approximation for the delay probability in the M/M/N queue, N large, which operates in the rationalized regime. (Some useful properties of $P(\cdot)$ and other functions are verified in Appendices A–C.) In §§6–8, we analyze the rationalized, efficiency- and quality-driven regimes, under the optimization approach. While the analysis is abstract, each of these sections concludes with examples of specific cost structures, for concreteness. In §9, we introduce the constraint satisfaction approach, which gives rise to the same three regimes of operation as optimization.

Section 10 describes numerical experiments that test the accuracy of our asymptotically supported approximations. As already mentioned, the findings are astounding—rarely do we miss by more than a single agent, as far as optimal staffing levels are concerned. In addition, even though the theory is asymptotic, our approximations are accurate with as few as three agents. In order to apply our approximations, guidelines are required for fitting a given call center, represented by its parameters and costs, to one of the three operational regimes. This turns out simpler than expected. Indeed, our numerical experiments, backed up by some theory, clearly establish the robustness of the rationalized approximation, as it covers accurately both the efficiency- and quality-driven regimes. Thus, except for extreme settings, the rationalized approximation is the one to use, as we do in the following section. We conclude in §11 with a few worthy directions for future research.

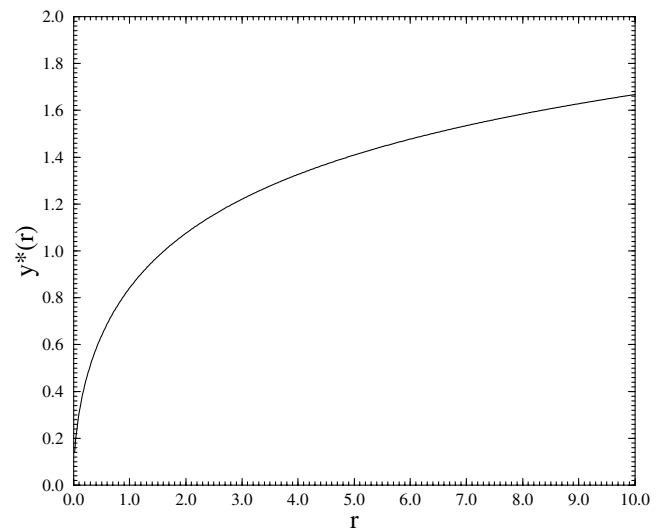
2. The Square-Root Safety Staffing Principle

To recapitulate, we determine asymptotically optimal staffing levels in accordance with the relative importance of agents’ costs and efficiency versus customers’ service quality. The very special case of linear staffing and delay costs (Example 6.3) already leads to the (re)discovery, as well as a deeper understanding, of a remarkably robust rule of thumb, the *square-root safety staffing* rule. It reads as follows: Suppose that the arrival rate is λ customers per hour, and service rate is μ , which implies that the system’s *offered load* is given by $R = \lambda/\mu$; if the staffing cost is $\$c$ per agent per hour, and waiting cost is $\$a$ per customer per hour, our recommended number of servers N^* is given by

$$N^* = R + y^*\left(\frac{a}{c}\right)\sqrt{R}, \quad (1)$$

where the function $y^*(r)$, $r \geq 0$, is plotted in Figures 1 and 2; see (23). In simple words, at least R agents ($\lfloor R \rfloor +$

Figure 1. $y^*(r)$ as function of r , $0 \leq r \leq 10$.

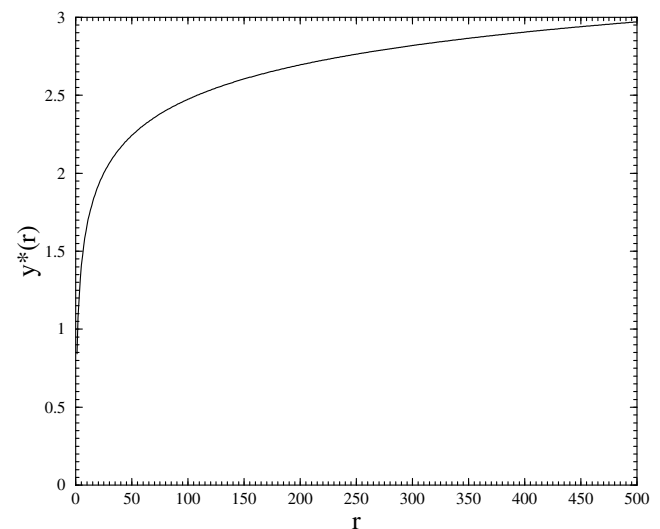


1 to be exact) are required to guarantee stability; however, safety staffing must be added to the minimum as a protection against stochastic variability. This number of additional agents is proportional to \sqrt{R} , and the proportionality coefficient $y^*(a/c)$ is determined through the optimization (23), by the relative importance of customers’ delay (a) to agents’ salary (c).

Note that the right-hand side of (1) need not be an integer, in which case N^* is obtained by rounding it off. We demonstrate in §10, below (40), that this yields the staffing level that minimizes waiting plus staffing costs, exactly in most cases, and off by a single agent in the other ones.

The *form* of (1) already carries with it important insight. Let $\Delta = y^*(a/c)\sqrt{R}$ denote the *safety staffing* level (the excess number of servers above the minimum $R = \lambda/\mu$). Then, with a and c fixed, an n -fold increase in the offered

Figure 2. $y^*(r)$ as function of r , $0 \leq r \leq 500$.



load R requires that the safety staffing Δ increases by only \sqrt{n} -fold, which constitutes significant economies of scale (Whitt 1992).

Now, suppose that R is measured in 100s, as it is in large call centers. Then \sqrt{R} is in the low 10s, hence Δ is as well (since y^* grows so slowly: $y^*(100) \approx 2.5$). It follows that the bulk of the agents, namely R , must be present for stability, and only a small fraction $\Delta/R = y^*/\sqrt{R}$ of these must be added as safety against stochastic variability (up to 10%, and, in fact, significantly less for large call centers, as the practical values of Δ/R indicate). This results in high agent utilization levels R/N^* —around 90% and up. Nevertheless, as shown in Examples 2.1–2.3, operational service quality ranges from the acceptable to the extremely high. (Indeed, small changes in Δ , which amounts to small changes in agents' utilization, have noticeable effects on performance.) Thus, we are operating in a regime where high resource utilization and service level coexist, which is due to economies of scale that dominate stochastic variability. It is important and interesting to note that data from large call centers confirms these observations (Garnett et al. 2002).

Small values of r correspond to efficiency-driven staffing. In this range, the function $y^*(\cdot)$ is reasonably approximated by

$$y^*(r) \approx \sqrt{\frac{r}{1 + r(\sqrt{\pi/2} - 1)}}, \quad 0 \leq r < 10.$$

Large values of r correspond to quality-driven staffing. In this range, a close lower bound is $y^*(r) \approx \sqrt{s - \ln s}$, where $s = 2 \ln(r/\sqrt{\pi})$, $r \uparrow \infty$. (See Remark 6.4 for some details on these asymptotic expansions.)

Under our square-root safety staffing, it is anticipated that service level, as expressed by the industry-standard TSF, equals

$$\text{TSF} = \Pr\{\text{Wait} > T\} \approx P(y^*) \times e^{-Ty^* \sqrt{\lambda\mu}}; \quad y^* = y^*\left(\frac{a}{c}\right);$$

in which

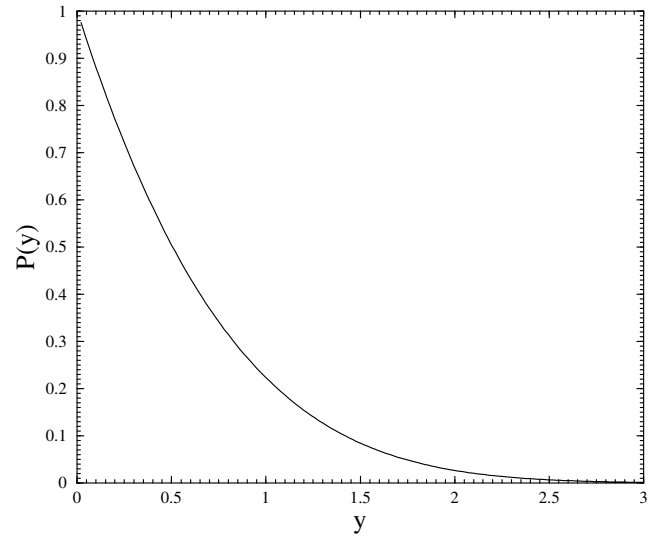
$$P(y^*) = \left[1 + \frac{y^* \Phi(y^*)}{\phi(y^*)}\right]^{-1} \approx \Pr\{\text{Wait} > 0\}, \quad (2)$$

is the *Halfin-Whitt delay function* (Halfin and Whitt 1981) (see Figure 3 and §5); $\phi(\cdot)$ and $\Phi(\cdot)$ are the density and cumulative distribution function of the standard normal distribution, respectively. A more management-friendly representation of TSF is

$$\begin{aligned} \text{TSF} &= \Pr\{\text{Wait} > T \times E[\text{Service Time}]\} \\ &\approx \Pr\{\text{Wait} > 0\} \times e^{-T\Delta}. \end{aligned} \quad (3)$$

Here delay is measured in units of average service time ($E[\text{Service Time}] = 1/\mu$), and $\Delta = y^* \sqrt{R}$ is the safety

Figure 3. The Halfin-Whitt delay function $P(y)$.



staffing level. Another service level standard is the average waiting time, often referred to as Average Speed of Answer (ASA). With N^* as in (1), and again naturally quantified in units of service durations, it is given by

$$\frac{\text{ASA}}{1/\mu} = \frac{E[\text{Wait}]}{E[\text{Service Time}]} \approx \frac{P(y^*)}{\Delta}. \quad (4)$$

The industry standard for measuring operational efficiency is agent utilization, namely R/N , which is traded off against service level. Agents are thus idle, or more appropriately described as being available for service, a fraction Δ/N of their time.

EXAMPLE 2.1. Consider, for example, the best-practice call center, described in the second paragraph of our Introduction. Assuming 1,800 calls per busy hour, the offered load equals $R = 120$. With 90% utilization, one expects that about $N^* \approx 133$ agents share the load ($\Delta = 13$), hence the center operates with $y^* \approx 1.22$. Inverting $y^*(\cdot)$ in Figure 1 shows that, in this call center, an hour wait of customers is valued as three times the hourly wage of an agent.

With this staffing level, it is expected that about 15% of the customers ($P(1.22) = 0.15$) are delayed; that 5% of the customers are delayed over 20 seconds (using (3) with $T = 1/12$); and that, by (4), ASA equals 2.7 seconds (while those who were delayed actually averaged 18 seconds waiting).

But the staffing level in the example can be interpreted differently. To this end, recall that the prevalent alternative to the above optimization approach is constraint satisfaction. Specifically, in Example 9.5 it is shown that the least N that guarantees $\Pr\{\text{Wait} > 0\} < \epsilon$ is closely approximated by rounding up

$$N^* = R + P^{-1}(\epsilon) \sqrt{R}, \quad (5)$$

where $P(\cdot)$ is the Halfin-Whitt delay function introduced in (2). Returning to the above best-practice call center, $P^{-1}(\epsilon) = 1.22$ yields, as expected, $\epsilon = 0.15$.

EXAMPLE 2.2. One should note that a constraint on the fraction of delayed customers is severe, hence it fits call centers that cater to, say, emergency calls. This can be nicely explained within our framework. For example, requiring that $\epsilon = 0.01$, namely one customer out of 100 delayed on average, corresponds to $y^* = P^{-1}(0.01) = 2.38$ (see Figure 3), which *could* be interpreted as saying (via Figure 2) that $a/c = (y^*)^{-1}(2.38) = 75$! An evaluation of customers' time as being worth 75-fold of agents' time seems reasonable only under extreme circumstances: For example, if the "servers" are "cheap" being, say, Interactive Voice Response (IVR) units, or customers' time is highly valued as with emergency call centers.

EXAMPLE 2.3. Most call centers define TSF with a positive T , and then requiring $\epsilon = 0.01$ need not be extreme. We now illustrate this by analyzing a prevalent industry standard, which is to aspire that no more than 80% of the callers are delayed over $T = 20$ seconds. Incidentally, we believe that the source of this standard is the familiar 20 : 80 managerial rule of thumb, stating in great generality and vagueness that "only 20% of the reasons already give rise to 80% of the problems." While there is no apparent reason for connecting this rule of thumb with any staffing standard, it is nevertheless worthwhile to note that our framework provides some interesting implications for using this rule. This will now be demonstrated via four scenarios which, for convenience, are also summarized in Table 1.

Consider a large call center with $\lambda = 100$ calls per minute, and 4 minutes average call duration. Thus $R = 400$, and adhering to the 20 : 80 rule implies that $y^* = 0.53$, hence, $N^* = 411$. By Figure 1, this translates into $a/c = 0.32$. It follows that, while customers are not highly valued, the 20 : 80 rule is "easy" to adhere to because of the call center's size. To wit, increasing N^* to 429 amounts to $y^* = 1.4$, or $a/c = 4.9$, reflecting a significant yet reasonable increase of the relative value of customers' to agents' time. This is accompanied by an increase in server availability (idleness), from 3% to 7%, which enables an order-of-magnitude reduction in TSF, from 0.2 to little less than 0.01: About one out of 100 customers is delayed for more than 20 seconds.

To underscore the role of scale in the above scenario, consider a call center with the same offered load parameters as Example 2.1: Thirty calls per minute, and again four minutes average call duration. Now $R = 120$, but it takes $N^* = 140$ to achieve $\text{TSF} = 0.01$, with $T = 20$ seconds. This corresponds to $y^* = 1.75$, or $a/c = 12.5$, a 2.5-fold increase over the large call center. It is interesting to note that with an average call duration of 30 seconds (as in 411 services), with T held at 20 seconds, $N^* = 126$ would suffice, which amounts to $y^* = 0.53$ and $a/c = 0.32$. This is identical to the large call center with the 20 : 80 rule operation, but the latter accommodates mean service time of four minutes, in contrast to the 30 seconds here.

The square-root safety staffing principle emerged from the simplest cost structure (linear staffing and waiting costs). While our framework accommodates general costs, the corresponding safety staffing levels are nevertheless *always* proportional to \sqrt{R} ; it is only the proportionality coefficient that varies with the cost.

2.1. Related Literature

The square-root safety staffing principle has been part of the queueing-theory folklore for a long time. Its origin goes back to Erlang's 1923 paper, published in Erlang (1948). Erlang derived the square-root principle via marginal analysis of the benefit in adding a server, indicating that it had been practiced, in fact, since 1913. More recently, the principle was well documented by Grassmann (1986, 1988) and then revisited by Kolesar and Green (1998), where both its accuracy and applicability have been convincingly confirmed. The principle was substantiated by Whitt (1992), then adapted in Jennings et al. (1996) to nonstationary models. Except for Erlang (1948), all the work we are aware of has applied infinite-server heuristics; it is grounded in the fact that the steady-state number of customers in the M/M/ ∞ queue, say Q^∞ , is Poisson distributed with mean $R = \lambda/\mu$. It follows that Q^∞ is approximately normally distributed, with mean R and standard deviation \sqrt{R} , when R is not too small. To relate this to staffing in the M/M/N model, one approximates the latter's probability of delay by

$$\Pr\{Q^\infty \geq N\} \approx 1 - \Phi\left(\frac{N - R}{\sqrt{R}}\right).$$

Table 1. Example 2.3—summary of scenarios.

TSF with $T = 20$ seconds							
λ (minute ⁻¹)	$1/\mu$ (minute)	$R = \lambda/\mu$	N^*	y^*	$\lambda/(N^*\mu)$	a/c	TSF
100	4	400	411	0.53	0.97	0.32	0.20
100	4	400	429	1.4	0.93	4.9	0.01
30	4	120	140	1.75	0.86	12.5	0.01
240	0.5	120	126	0.53	0.95	0.32	0.01

Then, the staffing level N^* that guarantees ϵ delay probability is chosen to be

$$N^* = R + \bar{\Phi}^{-1}(\epsilon)\sqrt{R}, \quad (6)$$

where $\bar{\Phi} = 1 - \Phi$.

The square-root principle contains two parts: First, the conceptual observation that the safety staffing level is proportional to the square-root of the offered load; and second, the explicit calculation of the proportionality coefficient y^* . Our framework accommodates *both* of these two needs, while in all previous works, to the best of our understanding, at least one of them is treated in a heuristic fashion or simply ignored. (We shall be specific shortly.) More important, however, is the fact that our approach and framework allow an arbitrary cost structure, and they have the potential to generalize beyond Erlang-C. For a concrete example, Garnett et al. (2002) accommodate impatient customers: In their main result, the square-root rule arises conceptually, but the determination of the value of y^* is left open. Being specific now, Whitt (1992) and Jennings et al. (1996) refer to y^* as a measure of service level, but leave out any explicit calculation of it. Grassmann (1988), taking the optimization approach, leads the reader through an instructive progression of increasingly complex staffing models, culminating in his “equilibrium model” (Erlang-C), for which no “square-root” justification is provided. (It is justified for his less complex model, under the “Independence Assumption,” but this amounts to using (6).) Some numerical experiments, inspired by Grassman (1988), are reported at the beginning of §10. Finally Kolesar and Green (1998) advocate the use of (6), in order to support constraint satisfaction that achieves $\Pr\{\text{Wait} > 0\} \leq \epsilon$. We, on the other hand, recommend the use of (5) for constraint satisfaction, which is proven asymptotically accurate in Example 9.5. The approximations (5) and (6) essentially coincide for small values of ϵ , but (5) is uniformly more accurate. We refer to the beginning of §10 for more details.

3. Model Description

We consider the classical M/M/N (Erlang-C) model with N servers and infinite-capacity waiting room. Customers arrive as a Poisson process of rate λ , and have independent exponentially distributed service times with mean $1/\mu$. The service rate μ will be arbitrary but fixed, whereas the arrival rate λ will grow large in order to obtain asymptotic scaling results. We assume $\lambda/N\mu < 1$ for stability. Customers are served in order of arrival; then (see, for instance, Cooper 1981) the waiting-time distribution is given by

$$\Pr\{\text{Wait} > t\} = \pi(N, \lambda/\mu) e^{-(N\mu - \lambda)t},$$

where the probability of waiting $\pi(N, \lambda/\mu) = \Pr\{\text{Wait} > 0\}$ is determined by

$$\pi(N, \nu) = \frac{\nu^N}{N!} \left\{ (1 - \nu/N) \sum_{n=0}^{N-1} \frac{\nu^n}{n!} + \frac{\nu^N}{N!} \right\}^{-1}.$$

We consider the problem of determining the staffing level N that optimally balances staffing cost against quality-of-service. To this end, a staffing cost $F(N)$ per unit of time is associated with staffing N servers. As mentioned in the Introduction, we assume that $F(N)$ is also defined for all noninteger values $N > \lambda/\mu$, and that this extended function $F(\cdot)$ is convex and strictly increasing, which also covers linear costs.

Quality-of-service is quantified in terms of a waiting-cost function $D_\lambda(\cdot)$: A cost $D_\lambda(t)$ is incurred when a customer waits for t time units. (The subscript λ is attached to allow for the possibility that the primitives vary with the arrival intensity.) We assume that $D_\lambda(\cdot)$ is strictly increasing. Without loss of generality, we may take $D_\lambda(0) = 0$. The expected total cost per unit of time is then given by

$$\begin{aligned} C(N, \lambda) &= F(N) + \lambda E[D_\lambda(\text{Wait})] \\ &= F(N) + \lambda \pi(N, \lambda/\mu) G(N, \lambda), \end{aligned} \quad (7)$$

where

$$\begin{aligned} G(N, \lambda) &= E[D_\lambda(\text{Wait}) \mid \text{Wait} > 0] \\ &= (N\mu - \lambda) \int_0^\infty D_\lambda(t) e^{-(N\mu - \lambda)t} dt. \end{aligned}$$

Notice that $G(N, \lambda)$ is also defined for all noninteger values $N > \lambda/\mu$. We assume that $D_\lambda(\cdot)$ is such that $G(N, \lambda)$ is finite for all $\lambda/\mu < N$.

We are interested in determining the optimum staffing level

$$N_\lambda^* := \arg \min_{N > \lambda/\mu} C(N, \lambda) \quad (8)$$

(the minimization being over integer values). To see that N_λ^* is well defined, notice that $\lim_{N \rightarrow \infty} F(N) = \infty$, and thus $\lim_{N \rightarrow \infty} C(N, \lambda) = \infty$. Hence, $C(N, \lambda)$ indeed achieves a minimum value.

4. Framework for Asymptotic Optimality

In principle, the optimum staffing level N_λ^* in Equation (8) may be obtained through brute-force enumeration. Rather than determining the optimum staffing level numerically, however, we are primarily interested in gaining insight into how N_λ^* grows with the arrival intensity λ , and how it depends on the staffing and waiting cost functions $F(N)$ and $G(N, \lambda)$. In order to do so, we develop an approximate analytical approach for determining the optimum staffing level. As a first step, we translate the discrete optimization problem (8) into a continuous one. The next step is to approximate the latter problem by a related continuous version, which is easier to solve. To validate the approach, we then prove that the optimal solution to the approximating continuous problem provides an asymptotically optimal solution to the original discrete problem.

We first transform the discrete optimization problem into a continuous one. Let

$$N_\lambda(x) = \lambda/\mu + x\sqrt{\lambda/\mu},$$

so that the variable $x = (N - \lambda/\mu)/\sqrt{\lambda/\mu}$ is the (normalized) number of servers in excess of the minimum number λ/μ required for stability. In terms of x , we define

$$F_\lambda(x) := F(N_\lambda(x)) - F(\lambda/\mu);$$

$$G_\lambda(x) := \lambda G(N_\lambda(x), \lambda);$$

$$C_\lambda(x) := C(N_\lambda(x), \lambda) - F(\lambda/\mu);$$

$$\pi_\lambda(x) := H(N_\lambda(x), \lambda/\mu),$$

with

$$H(M, \alpha) = \left\{ \alpha \int_0^\infty e^{-\alpha t} t(1+t)^{M-1} dt \right\}^{-1}.$$

It can be verified (Jagers and Van Doorn 1986, 1991) that $\pi_\lambda(x) = \pi(N_\lambda(x), \lambda/\mu)$ for integer values of $N_\lambda(x)$. The total cost per unit of time (up to the additive constant factor $F(\lambda/\mu)$) can thus be rewritten

$$C_\lambda(x) = F_\lambda(x) + \pi_\lambda(x)G_\lambda(x).$$

Denote

$$x_\lambda^* := \arg \min_{x>0} C_\lambda(x). \quad (9)$$

To see that x_λ^* is well defined, first notice that the function $C_\lambda(\cdot)$ is strictly convex. This follows from the assumption that $F(\cdot)$ is convex and the fact that $\pi_\lambda(\cdot)$ is convex (Jagers and Van Doorn 1986, 1991) and $G_\lambda(\cdot)$ is strictly convex (Appendix C). In addition, $\lim_{x \downarrow 0} C_\lambda(x) = \infty$, since $\lim_{N \downarrow \lambda/\mu} G(N, \lambda) = \infty$. Also, $\lim_{x \rightarrow \infty} C_\lambda(x) = \infty$, because $\lim_{N \rightarrow \infty} F(N) = \infty$. Hence, $C_\lambda(\cdot)$ is unimodal, implying that it, indeed, achieves a unique minimum value at $x_\lambda^* \in (0, \infty)$. Further, notice that either $N_\lambda^* = \lfloor N_\lambda(x_\lambda^*) \rfloor$ or $N_\lambda^* = \lceil N_\lambda(x_\lambda^*) \rceil$, which establishes the link between the discrete problem and the corresponding continuous problem. (Here $\lfloor u \rfloor$ and $\lceil u \rceil$ denote the largest integer smaller than or equal to u , and the smallest integer larger than or equal to u , respectively.) Next, we approximate x_λ^* in (9) by

$$z_\lambda^* := \arg \min_{z>0} C[z; \hat{F}_\lambda, \hat{\pi}_\lambda, \hat{G}_\lambda], \quad (10)$$

where

$$C[z; \hat{F}_\lambda, \hat{\pi}_\lambda, \hat{G}_\lambda] := \hat{F}_\lambda(z) + \hat{\pi}_\lambda(z)\hat{G}_\lambda(z),$$

with the functions $\hat{F}_\lambda(\cdot)$, $\hat{\pi}_\lambda(\cdot)$, $\hat{G}_\lambda(\cdot)$ “approximating” $F_\lambda(\cdot)$, $\pi_\lambda(\cdot)$, $G_\lambda(\cdot)$, respectively. (Note that with this notation, $x_\lambda^* = \arg \min_{x>0} C[x; F_\lambda, \pi_\lambda, G_\lambda]$.) The approximating

functions $\hat{F}_\lambda(\cdot)$, $\hat{G}_\lambda(\cdot)$, and $\hat{\pi}_\lambda(\cdot)$ that we consider will always be such that z_λ^* exists and is unique. If $\hat{F}_\lambda(\cdot)$, $\hat{G}_\lambda(\cdot)$, $\hat{\pi}_\lambda(\cdot)$ have a simple form, then solving for z_λ^* will be easier than determining x_λ^* . At the same time, if $\hat{F}_\lambda(\cdot)$, $\hat{G}_\lambda(\cdot)$, and $\hat{\pi}_\lambda(\cdot)$ approximate $F_\lambda(\cdot)$, $G_\lambda(\cdot)$, and $\pi_\lambda(\cdot)$ well, then it is reasonable to expect that z_λ^* provides a good approximation to x_λ^* and, moreover, $N_\lambda(z_\lambda^*)$ yields a good approximation to N_λ^* .

Before formalizing the above approximation principle, we introduce the following notational conventions: For any pair of functions a_λ and b_λ (implicitly assuming existence of the limits), denote

$$a_\lambda \approx b_\lambda: \lim_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} = 1; \quad a_\lambda \approx \gamma b_\lambda: \lim_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} = \gamma, \quad 0 < \gamma < \infty;$$

$$a_\lambda \ll b_\lambda: \lim_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} = 0; \quad a_\lambda \gg b_\lambda: \lim_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} = \infty;$$

$$a_\lambda \leq b_\lambda: \limsup_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} \leq 1; \quad a_\lambda \geq b_\lambda: \liminf_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} \geq 1;$$

$$a_\lambda < b_\lambda: \liminf_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} < 1; \quad a_\lambda > b_\lambda: \limsup_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} > 1;$$

$$a_\lambda \ll b_\lambda: \liminf_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} = 0; \quad a_\lambda \gg b_\lambda: \limsup_{\lambda \rightarrow \infty} \frac{a_\lambda}{b_\lambda} = \infty.$$

LEMMA 4.1. Denote $\hat{C}_\lambda(z) = C[z; \hat{F}_\lambda, \hat{\pi}_\lambda, \hat{G}_\lambda]$. Then $C_\lambda(z_\lambda^*) \approx C_\lambda(x_\lambda^*)$ if both $C_\lambda(x_\lambda^*) \approx \hat{C}_\lambda(x_\lambda^*)$ and $C_\lambda(z_\lambda^*) \approx \hat{C}_\lambda(z_\lambda^*)$.

PROOF. By definition of x_λ^* , $C_\lambda(z_\lambda^*) \geq C_\lambda(x_\lambda^*)$, so it suffices to show that $C_\lambda(z_\lambda^*) \leq C_\lambda(x_\lambda^*)$, which follows directly from

$$C_\lambda(z_\lambda^*) \approx \hat{C}_\lambda(z_\lambda^*) \leq \hat{C}_\lambda(x_\lambda^*) \approx C_\lambda(x_\lambda^*). \quad \square$$

Define

$$S_\lambda(x) := \min\{C(\lfloor N_\lambda(x) \rfloor, \lambda), C(\lceil N_\lambda(x) \rceil, \lambda)\}. \quad (11)$$

LEMMA 4.2. If $C_\lambda(z_\lambda^*) \approx C_\lambda(x_\lambda^*)$, then $S_\lambda(z_\lambda^*) - F(\lambda/\mu) \approx C(N_\lambda^*, \lambda) - F(\lambda/\mu)$.

PROOF. By definition, $S_\lambda(z_\lambda^*) \geq C(N_\lambda^*, \lambda)$, so it suffices to show that

$$S_\lambda(z_\lambda^*) - F(\lambda/\mu) \leq C(N_\lambda^*, \lambda) - F(\lambda/\mu).$$

For fixed λ , we distinguish between four cases.

(i) $N_\lambda^* - 1 < N_\lambda(z_\lambda^*) \leq N_\lambda^*$. Then $\lfloor N_\lambda(z_\lambda^*) \rfloor = N_\lambda^*$, and $S_\lambda(z_\lambda^*) = C(N_\lambda^*, \lambda)$.

(ii) $N_\lambda^* \leq N_\lambda(z_\lambda^*) < N_\lambda^* + 1$. Then $\lceil N_\lambda(z_\lambda^*) \rceil = N_\lambda^*$, and $S_\lambda(z_\lambda^*) = C(N_\lambda^*, \lambda)$.

(iii) $N_\lambda(z_\lambda^*) \leq N_\lambda^* - 1$. Then

$$\begin{aligned} z_\lambda^* &\leq \frac{\lceil N_\lambda(z_\lambda^*) \rceil - \lambda/\mu}{\sqrt{\lambda/\mu}} \leq \frac{N_\lambda^* - 1 - \lambda/\mu}{\sqrt{\lambda/\mu}} \\ &\leq \frac{\lfloor N_\lambda(x_\lambda^*) \rfloor - \lambda/\mu}{\sqrt{\lambda/\mu}} \leq x_\lambda^*, \end{aligned}$$

so that

$$\begin{aligned} S_\lambda(z_\lambda^*) - F(\lambda/\mu) &\leq C(\lceil N_\lambda(z_\lambda^*) \rceil, \lambda) - F(\lambda/\mu) \\ &= C_\lambda\left(\frac{\lceil N_\lambda(z_\lambda^*) \rceil - \lambda/\mu}{\sqrt{\lambda/\mu}}\right) \leq C_\lambda(z_\lambda^*) \end{aligned}$$

because of the unimodality of $C_\lambda(\cdot)$.

(iv) $N_\lambda(z_\lambda^*) \geq N_\lambda^* + 1$. Then

$$\begin{aligned} z_\lambda^* &\geq \frac{\lfloor N_\lambda(z_\lambda^*) \rfloor - \lambda/\mu}{\sqrt{\lambda/\mu}} \geq \frac{N_\lambda^* + 1 - \lambda/\mu}{\sqrt{\lambda/\mu}} \\ &\leq \frac{\lfloor N_\lambda(x_\lambda^*) \rfloor - \lambda/\mu}{\sqrt{\lambda/\mu}} \leq x_\lambda^*, \end{aligned}$$

so that

$$\begin{aligned} S_\lambda(z_\lambda^*) - F(\lambda/\mu) &\leq C(\lfloor N_\lambda(z_\lambda^*) \rfloor, \lambda) - F(\lambda/\mu) \\ &= C_\lambda\left(\frac{\lfloor N_\lambda(z_\lambda^*) \rfloor - \lambda/\mu}{\sqrt{\lambda/\mu}}\right) \leq C_\lambda(z_\lambda^*). \end{aligned}$$

Thus, for all λ ,

$$\begin{aligned} S_\lambda(z_\lambda^*) - F(\lambda/\mu) &\leq \max\{C(N_\lambda^*, \lambda) - F(\lambda/\mu), C_\lambda(z_\lambda^*)\} \\ &\approx \max\{C(N_\lambda^*, \lambda) - F(\lambda/\mu), C_\lambda(x_\lambda^*)\} \\ &= C(N_\lambda^*, \lambda) - F(\lambda/\mu). \quad \square \end{aligned}$$

Combining Lemmas 4.1 and 4.2, we obtain the fundamental approximation principle underlying our approach:

COROLLARY 4.3 (ASYMPTOTIC OPTIMALITY). Denote $\hat{C}_\lambda(z) = C[z; \hat{F}_\lambda, \hat{\pi}_\lambda, \hat{G}_\lambda]$. Let x_λ^* and z_λ^* be as in (9) and (10), respectively. If $C_\lambda(x_\lambda^*) \approx \hat{C}_\lambda(x_\lambda^*)$ and $C_\lambda(z_\lambda^*) \approx \hat{C}_\lambda(z_\lambda^*)$, then the staffing function z_λ^* is asymptotically optimal in the sense that, as $\lambda \rightarrow \infty$,

$$S_\lambda(z_\lambda^*) - F(\lambda/\mu) \approx C(N_\lambda^*, \lambda) - F(\lambda/\mu),$$

with $S_\lambda(x)$ given in (11).

Note that the quantities $S_\lambda(z_\lambda^*) - F(\lambda/\mu)$ and $C(N_\lambda^*, \lambda) - F(\lambda/\mu)$ may be interpreted as the total cost in excess of the minimum required staffing cost $F(\lambda/\mu)$ for the approximately optimal staffing level $N_\lambda(z_\lambda^*)$ and for the truly optimal level N_λ^* , respectively. The above corollary identifies conditions under which these two quantities are asymptotically equal, implying that the approximate solution is asymptotically optimal in a certain sense.

In the next sections, we identify “simple” functions $\hat{F}_\lambda(\cdot)$, $\hat{G}_\lambda(\cdot)$, and $\hat{\pi}_\lambda(\cdot)$, such that $F_\lambda(x_\lambda^*) \approx \hat{F}_\lambda(x_\lambda^*)$ and $F_\lambda(z_\lambda^*) \approx \hat{F}_\lambda(z_\lambda^*)$, $G_\lambda(x_\lambda^*) \approx \hat{G}_\lambda(x_\lambda^*)$ and $G_\lambda(z_\lambda^*) \approx \hat{G}_\lambda(z_\lambda^*)$, $\pi_\lambda(x_\lambda^*) \approx \hat{\pi}_\lambda(x_\lambda^*)$ and $\pi_\lambda(z_\lambda^*) \approx \hat{\pi}_\lambda(z_\lambda^*)$. This implies $C_\lambda(x_\lambda^*) \approx \hat{C}_\lambda(x_\lambda^*)$ and $C_\lambda(z_\lambda^*) \approx \hat{C}_\lambda(z_\lambda^*)$ as required in the above corollary, which then enables us to gain insight into the behavior of N_λ^* , as a function of λ .

5. Some Special Functions

In this section, we introduce some functions that will play a central role in our analysis.

For any $x > 0$, define

$$P(x) := \frac{1}{1 + (x/h(-x))}, \quad (12)$$

where $h(\cdot)$ is the “hazard rate” function of the standard normal distribution, namely

$$h(x) := \frac{\phi(x)}{1 - \Phi(x)},$$

with

$$\phi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}; \quad \Phi(x) = \int_{-\infty}^x \phi(y) dy.$$

In Lemma B.1 we prove that $P(\cdot)$ is strictly convex decreasing.

Also define

$$Q_\lambda(x) := \frac{\exp\{N_\lambda(x)[1 - r_\lambda(x) + \log r_\lambda(x)]\}}{\sqrt{2\pi N_\lambda(x)}(1 - r_\lambda(x))},$$

with

$$r_\lambda(x) := \frac{\lambda/\mu}{N_\lambda(x)},$$

and let

$$Q(x) := \frac{\phi(x)}{x} = \frac{e^{-x^2/2}}{x\sqrt{2\pi}}. \quad (13)$$

The following two lemmas characterize the asymptotic behavior of $\pi_\lambda(\cdot)$, as $\lambda \rightarrow \infty$.

LEMMA 5.1 (HALFIN AND WHITT 1981). For any function x_λ with $\limsup_{\lambda \rightarrow \infty} x_\lambda < \infty$,

$$\pi_\lambda(x_\lambda) \approx P(x_\lambda).$$

If, moreover, $\lim_{\lambda \rightarrow \infty} x_\lambda = x$, then $\pi_\lambda(x_\lambda) \approx P(x)$, $x \geq 0$. In particular, if $\lim_{\lambda \rightarrow \infty} x_\lambda = 0$, then $\pi_\lambda(x_\lambda) \approx 1$.

PROOF. Suppose to the contrary. Then there must be a subsequence $\{\lambda_n\}$ with $\lim_{n \rightarrow \infty} \lambda_n = \infty$ such that $\lim_{n \rightarrow \infty} x_{\lambda_n} = \beta$ and $\lim_{n \rightarrow \infty} \pi_{\lambda_n}(x_{\lambda_n}) = \alpha$, where $0 < \beta < \infty$ and $\alpha \neq P(\beta)$. This is in contradiction with Proposition 1 of Halfin and Whitt (1981), which asserts that $\alpha = P(\beta)$ must prevail for such a sequence $\{\lambda_n\}$. \square

LEMMA 5.2 (APPENDIX A). For any function x_λ with $\lim_{\lambda \rightarrow \infty} x_\lambda = \infty$,

$$\pi_\lambda(x_\lambda) \approx Q_\lambda(x_\lambda).$$

If also $x_\lambda \leq \lambda^{1/6}$, then

$$\pi_\lambda(x_\lambda) \approx Q(x_\lambda).$$

If specifically $x_\lambda = \kappa\sqrt{\lambda/\mu}$ for some constant $\kappa > 0$, then

$$\pi_\lambda(x_\lambda) \approx \frac{1}{\kappa\sqrt{2\pi\lambda/\mu}(1 + \kappa)} \left(\frac{e^\kappa}{(1 + \kappa)^{1+\kappa}} \right)^{\lambda/\mu}.$$

We conclude the section with some observations on the behavior of the functions $F_\lambda(\cdot)$, $G_\lambda(\cdot)$, and $\pi_\lambda(\cdot)$, as defined at the outset of §4.

Recall that the staffing cost function $F(\cdot)$ is convex increasing, which implies that the function $F_\lambda(\cdot)$ is convex increasing as well. In addition, $F_\lambda(0) = 0$. Hence, $F_\lambda(x)/F_\lambda(y) \leq x/y$ for any pair of numbers $x \leq y$. Thus,

$$a_\lambda \sup > b_\lambda \implies F_\lambda(a_\lambda) \sup > F_\lambda(b_\lambda), \quad (14)$$

and

$$a_\lambda \sup \gg b_\lambda \implies F_\lambda(a_\lambda) \sup \gg F_\lambda(b_\lambda). \quad (15)$$

Also, from Lemmas 5.1 and B.1, for fixed $b \geq 0$,

$$a_\lambda \inf < b \implies P(a_\lambda) \sup > P(b), \quad \pi_\lambda(a_\lambda) \sup > \pi_\lambda(b), \quad (16)$$

$$a_\lambda \sup > b \implies P(a_\lambda) \inf < P(b), \quad \pi_\lambda(a_\lambda) \inf < \pi_\lambda(b), \quad (17)$$

and noting that $\lim_{x \rightarrow \infty} \pi_\lambda(x) = 0$,

$$a_\lambda \sup \gg b \implies P(a_\lambda) \inf \ll P(b), \quad \pi_\lambda(a_\lambda) \inf \ll \pi_\lambda(b). \quad (18)$$

6. Case I: Rationalized Regime

In this section, we consider what we call a *rationalized* scenario, by which we mean that, for some $\kappa > 0$,

$$F_\lambda(\kappa) \approx G_\lambda(\kappa), \quad (19)$$

or in words, the staffing cost $F_\lambda(\cdot)$ is comparable to the waiting cost $G_\lambda(\cdot)$, as $\lambda \rightarrow \infty$.

For any $\lambda > 0$, define

$$y_\lambda^* := \arg \min_{y>0} C[y; F_\lambda, P, G_\lambda], \quad (20)$$

with $P(\cdot)$ as in (12).

THEOREM 6.1. *The staffing function y_λ^* is asymptotically optimal in the sense of Corollary 4.3.*

PROOF. The idea of the proof may be described as follows. In order for Corollary 4.3 to apply, we need to show that the function $P(\cdot)$ is an asymptotically exact approximation for the function $\pi_\lambda(\cdot)$ around the points x_λ^* and y_λ^* . In view of Lemma 5.1, it suffices to show that $\limsup_{\lambda \rightarrow \infty} x_\lambda^* < \infty$ and $\limsup_{\lambda \rightarrow \infty} y_\lambda^* < \infty$. We prove this by contradiction, arguing that if this were not the case, then just the staffing cost alone would be overwhelmingly larger than the total cost associated with the fixed staffing function κ , contradicting the supposed optimality of x_λ^* or y_λ^* .

We start with showing that $\limsup_{\lambda \rightarrow \infty} x_\lambda^* < \infty$. Suppose to the contrary. Then $x_\lambda^* \sup \gg \kappa$, so that

$$F_\lambda(x_\lambda^*) \sup \gg F_\lambda(\kappa)$$

from (15). Using (19),

$$F_\lambda(\kappa) \approx F_\lambda(\kappa) + G_\lambda(\kappa) \geq F_\lambda(\kappa) + \pi_\lambda(\kappa)G_\lambda(\kappa) = C_\lambda(\kappa).$$

By definition,

$$C_\lambda(x_\lambda^*) \geq F_\lambda(x_\lambda^*).$$

Combining the above relations, we deduce

$$C_\lambda(x_\lambda^*) \sup \gg C_\lambda(\kappa),$$

contradicting the optimality of x_λ^* . Thus, $\limsup_{\lambda \rightarrow \infty} x_\lambda^* < \infty$. By similar arguments, $\limsup_{\lambda \rightarrow \infty} y_\lambda^* < \infty$.

Hence, according to Lemma 5.1, $\pi_\lambda(x_\lambda^*) \approx \hat{\pi}_\lambda(x_\lambda^*) = P(x_\lambda^*)$ and $\pi_\lambda(y_\lambda^*) \approx \hat{\pi}_\lambda(y_\lambda^*) = P(y_\lambda^*)$. Applying Corollary 4.3 then completes the proof. \square

PROPOSITION 6.2. *Assume that there exist functions $f(\cdot)$, $g(\cdot)$, and H_λ such that, for any function $k_\lambda > 0$,*

$$F_\lambda(k_\lambda) \approx f(k_\lambda)H_\lambda, \quad (21)$$

and

$$G_\lambda(k_\lambda) \approx g(k_\lambda)H_\lambda, \quad (22)$$

so certainly $F_\lambda(\kappa) \approx G_\lambda(\kappa)$ for all $\kappa > 0$.

Define

$$y^* = \arg \min_{y>0} C[y; f, P, g].$$

The staffing function y^* is then asymptotically optimal in the sense of Corollary 4.3.

PROOF. Similar to that of Theorem 6.1. Note that $F_\lambda(x_\lambda^*) \approx \hat{F}_\lambda(x_\lambda^*) = f(x_\lambda^*)H_\lambda$, $F_\lambda(y^*) \approx \hat{F}_\lambda(y^*) = f(y^*)H_\lambda$, $G_\lambda(x_\lambda^*) \approx \hat{G}_\lambda(x_\lambda^*) = g(x_\lambda^*)H_\lambda$, and $G_\lambda(y^*) \approx \hat{G}_\lambda(y^*) = g(y^*)H_\lambda$. \square

EXAMPLE 6.3. Assume there is a staffing cost c_λ per server per time unit, and a waiting cost a_λ per customer per time unit, as well as a fixed penalty cost b_λ when the waiting time exceeds d_λ time units, i.e., $F(N) := Nc_\lambda$ and $D_\lambda(t) = a_\lambda t + b_\lambda \mathbf{I}_{\{t > d_\lambda\}}$, so that $G(N, \lambda) = a_\lambda / (N\mu - \lambda) + b_\lambda e^{-(N\mu - \lambda)d_\lambda}$. Thus $F_\lambda(\kappa) = c_\lambda \kappa \sqrt{\lambda/\mu}$ and $G_\lambda(\kappa) = a_\lambda \sqrt{\lambda/\mu} / \kappa + b_\lambda \lambda e^{-\mu \kappa d_\lambda \sqrt{\lambda/\mu}}$.

(i) First, suppose that $a_\lambda \approx a$, $b_\lambda \sqrt{\lambda} e^{-\mu \kappa d_\lambda \sqrt{\lambda/\mu}} \ll 1$, and $c_\lambda \approx c$. Then (21)–(22) are satisfied for $f(\kappa) = c\kappa$, $g(\kappa) = a/\kappa$, and $H_\lambda = \sqrt{\lambda/\mu}$. Proposition 6.2 says that the staffing function

$$y^* = \arg \min_{y>0} \left\{ cy + \frac{aP(y)}{y} \right\} \quad (23)$$

is asymptotically optimal in the sense of Corollary 4.3. This is exactly the $y^*(a/c)$ that appears in (1). The numerical search for y^* is straightforward, since the function that it minimizes is unimodal. It is important enough for our purposes, as demonstrated in §2, that we plotted it in Figures 1 and 2.

(ii) Now, suppose that $a_\lambda \ll 1$, $b_\lambda \approx b/\sqrt{\lambda}$, $c_\lambda \approx c$, and $d_\lambda \approx d/\sqrt{\lambda}$ with $bd\mu > c$. Then (21)–(22) are satisfied for $f(\kappa) = c\kappa$, $g(\kappa) = b\sqrt{\mu}e^{-d\kappa\sqrt{\mu}}$, and $H_\lambda = \sqrt{\lambda/\mu}$. The asymptotically optimal staffing function is

$$y^* = \arg \min_{y>0} \{cy + bP(y)\sqrt{\mu}e^{-dy\sqrt{\mu}}\}.$$

(iii) Finally, consider the “combined” case where all costs show up in the limit, i.e., suppose that $a_\lambda \approx a$, $b_\lambda \approx b/\sqrt{\lambda}$, $c_\lambda \approx c$, and $d_\lambda \approx d/\sqrt{\lambda}$. Then (21)–(22) are satisfied for $f(\kappa) = c\kappa$, $g(\kappa) = (a/\kappa) + b\sqrt{\mu}e^{-d\kappa\sqrt{\mu}}$, and $H_\lambda = \sqrt{\lambda/\mu}$. The asymptotically optimal staffing function is

$$y^* = \arg \min_{y>0} \left\{ cy + P(y) \left[\frac{a}{y} + b\sqrt{\mu}e^{-dy\sqrt{\mu}} \right] \right\}.$$

REMARK 6.4 (ASYMPTOTIC EXPANSIONS FOR $y^*(\cdot)$). Two asymptotic expansions for $y^*(r)$ were quoted in §2, one for small r and the other for large. To derive the former, one simply replaces $P(y)$ in (23) by $P(0) + yP'(0) + (1/2)y^2P''(0)$, then uses the values for the derivatives from Appendix B, and finally minimizes the resulting simple function. (This yields $y^*(r)$, but with $\pi/2$ instead of the presently used $\sqrt{\pi/2}$. Empirically, the former was found to be more accurate near the origin, say $r \leq 0.5$, but the latter is a better approximation over the range $0 \leq r \leq 10$.)

As for large values, one uses the well-known approximation $1 - \Phi(y) \approx \phi(y)/y$ to get that also $P(y) \approx \phi(y)/y$. Substituting this approximation for $P(y)$ into (23) identifies y^* as the solution y of $ye^{y^2/2} = (a/\sqrt{2\pi})$.

Changing variables to $x = y^2$, then squaring, gives $x(t)e^{x(t)} = t$, where $t = (a^2/2\pi)$. A complete asymptotic expansion of $x(t)$, for large t , is calculated in De Bruijn (1981, pp. 25–28). This would yield the formula quoted in the Introduction, but with $s = \ln(r/\sqrt{2\pi})$. Again, it was found empirically that $s = \ln(r/\sqrt{\pi})$ provides a better approximation for $20 \leq r \leq 500$.

7. Case II: Efficiency-Driven Regime

In this section, we consider an *efficiency-driven* scenario, meaning that, for all $\kappa > 0$,

$$F_\lambda(\kappa) \gg G_\lambda(\kappa), \quad (24)$$

($\limsup_{\lambda \rightarrow \infty} (G_\lambda(\kappa)/F_\lambda(\kappa)) < \infty$ is actually sufficient) or in words, the staffing cost dominates the waiting cost, as $\lambda \rightarrow \infty$.

For any $\lambda > 0$, define

$$y_\lambda^* := \arg \min_{y>0} C[y; F_\lambda, 1, G_\lambda]. \quad (25)$$

THEOREM 7.1. *The staffing level y_λ^* is asymptotically optimal in the sense of Corollary 4.3.*

PROOF. The idea of the proof is largely similar to that of Theorem 6.1. We start with showing that $\lim_{\lambda \rightarrow \infty} x_\lambda^* = 0$. Suppose to the contrary. Then $x_\lambda^* \sup > u$ for some $u > 0$, so that $F_\lambda(x_\lambda^*) \sup > F_\lambda(u)$ from (14). Using (24), $F_\lambda(u) \approx F_\lambda(u) + G_\lambda(u) \geq F_\lambda(u) + \pi_\lambda(u)G_\lambda(u) = C_\lambda(u)$.

Combining the above relations, we obtain $C_\lambda(x_\lambda^*) \geq F_\lambda(x_\lambda^*) \sup > C_\lambda(u)$, contradicting the optimality of x_λ^* . Thus, $\lim_{\lambda \rightarrow \infty} x_\lambda^* = 0$. By similar arguments, $\lim_{\lambda \rightarrow \infty} y_\lambda^* = 0$. Hence, according to Lemma 5.1, $\pi_\lambda(x_\lambda^*) \approx \hat{\pi}_\lambda(x_\lambda^*) = P(x_\lambda^*) \approx 1$, and $\pi_\lambda(y_\lambda^*) \approx \hat{\pi}_\lambda(y_\lambda^*) = P(y_\lambda^*) \approx 1$.

Applying Corollary 4.3 then completes the proof. \square

PROPOSITION 7.2. *Assume that there exist functions $f(\cdot)$, $g(\cdot)$, H_λ , and J_λ such that, for any function $k_\lambda > 0$,*

$$F_\lambda(k_\lambda) \approx f(k_\lambda)H_\lambda, \quad (26)$$

and

$$G_\lambda(k_\lambda) \approx g(k_\lambda)H_\lambda J_\lambda, \quad (27)$$

with $J_\lambda \ll 1$, so certainly $F_\lambda(\kappa) \gg G_\lambda(\kappa)$ for all $\kappa > 0$.

Define $y_\lambda^* = \arg \min_{y>0} C[y; f, 1, gJ_\lambda]$. Assume that $f(\cdot)$ is convex increasing and that $g(\cdot)$ is strictly convex decreasing, with $\lim_{x \downarrow 0} g(x) = \infty$ so that y_λ^* exists and is unique.

The staffing function y_λ^* is then asymptotically optimal in the sense of Corollary 4.3.

PROOF. Similar to that of Theorem 7.1. Note that $F_\lambda(x_\lambda^*) \approx \hat{F}_\lambda(x_\lambda^*) = f(x_\lambda^*)H_\lambda$, $F_\lambda(y_\lambda^*) \approx \hat{F}_\lambda(y_\lambda^*) = f(y_\lambda^*)H_\lambda$, $G_\lambda(x_\lambda^*) \approx \hat{G}_\lambda(x_\lambda^*) = g(x_\lambda^*)H_\lambda J_\lambda$, and $G_\lambda(y_\lambda^*) \approx \hat{G}_\lambda(y_\lambda^*) = g(y_\lambda^*)H_\lambda J_\lambda$. \square

REMARK 7.3. The rationalized staffing function (20) is, in fact, also asymptotically optimal in the efficiency-driven regime as defined by (24). However, the staffing function (25), where $P(\cdot)$ is replaced by 1, is asymptotically optimal as well, while considerably simpler.

EXAMPLE 7.4. Consider the same cost structure as in Example 6.3. Now, however, suppose that $a_\lambda \approx aJ_\lambda$, $b_\lambda \approx bJ_\lambda/\sqrt{\lambda}$, $c_\lambda \approx c$, $d_\lambda \approx d/\sqrt{\lambda}$, and $J_\lambda \ll 1$. Then (26)–(27) are satisfied for $f(\kappa) = c\kappa$, $g(\kappa) = (a/\kappa) + b\sqrt{\mu}e^{-d\kappa\sqrt{\mu}}$, and $H_\lambda = \sqrt{\lambda/\mu}$. The asymptotically optimal staffing function is

$$y_\lambda^* = \arg \min_{y>0} \left\{ cy + \left[\frac{a}{y} + b\sqrt{\mu}e^{-dy\sqrt{\mu}} \right] J_\lambda \right\}.$$

In the special case where the nonlinear penalty cost is asymptotically negligible, i.e., “ $b = 0$ ” (or rather $b_\lambda\sqrt{\lambda}e^{-\mu\kappa d_\lambda\sqrt{\lambda/\mu}} \ll J_\lambda$), the optimal staffing function reduces to

$$y_\lambda^* = \arg \min_{y>0} \left\{ cy + \frac{a}{y} J_\lambda \right\}.$$

Note that the extreme case where the linear waiting cost is asymptotically vanishing, i.e., “ $a = 0$ ” (formally, $a_\lambda \ll J_\lambda$), does not make sense, since $\lim_{\kappa \downarrow 0} g(\kappa) = b\sqrt{\mu} < \infty$. Thus, the optimization problem will not have a strictly positive solution for large λ .

8. Case III: Quality-Driven Regime

In this section, we consider a *quality-driven* regime, meaning that, for all $\kappa > 0$,

$$F_\lambda(\kappa) \ll G_\lambda(\kappa), \quad (28)$$

or in words, the staffing cost is negligible compared to the waiting cost, as $\lambda \rightarrow \infty$.

For any $\lambda > 0$, define $y_\lambda^* := \arg \min_{y>0} C[y; F_\lambda, Q_\lambda, G_\lambda]$.

THEOREM 8.1. *The staffing function y_λ^* is asymptotically optimal in the sense of Corollary 4.3.*

PROOF. The idea of the proof is largely similar to that of Theorem 6.1. However, now we rely on Lemma 5.2, which indicates that the function $Q_\lambda(\cdot)$ is an asymptotically exact approximation for the function $\pi_\lambda(\cdot)$ around the points x_λ^* and y_λ^* , provided $\lim_{\lambda \rightarrow \infty} x_\lambda^* = \infty$ and $\lim_{\lambda \rightarrow \infty} y_\lambda^* = \infty$. This may be shown again by contradiction, verifying that if this were not the case then just the waiting cost alone would now be larger than the total cost associated with some fixed staffing function u .

We start with showing that $\lim_{\lambda \rightarrow \infty} x_\lambda^* = \infty$. Suppose to the contrary. Then $x_\lambda^* < u$ for some $u > 0$, so that $\pi_\lambda(x_\lambda^*)G_\lambda(x_\lambda^*) \stackrel{\sup}{>} \pi_\lambda(u)G_\lambda(u)$ from (16) and the fact that $G_\lambda(\cdot)$ is decreasing. Using Lemma 5.2 and (28),

$$\begin{aligned} \pi_\lambda(u)G_\lambda(u) &\stackrel{\sim}{\sim} P(u)G_\lambda(u) \stackrel{\sim}{\sim} F_\lambda(u) + P(u)G_\lambda(u) \\ &\stackrel{\sim}{\sim} F_\lambda(u) + \pi_\lambda(u)G_\lambda(u) = C_\lambda(u). \end{aligned}$$

Combining the above relations, we deduce that $C_\lambda(x_\lambda^*) \geq \pi_\lambda G_\lambda(x_\lambda^*) \stackrel{\sup}{>} C_\lambda(u)$, contradicting the optimality of x_λ^* . Thus, $\lim_{\lambda \rightarrow \infty} x_\lambda^* = \infty$. By similar arguments, $\lim_{\lambda \rightarrow \infty} y_\lambda^* = \infty$.

Hence, according to Lemma 5.2, $\pi_\lambda(x_\lambda^*) \stackrel{\sim}{\sim} \hat{\pi}_\lambda(y_\lambda^*) = Q_\lambda(x_\lambda^*)$, and $\pi_\lambda(y_\lambda^*) \stackrel{\sim}{\sim} \hat{\pi}_\lambda(y_\lambda^*) = Q_\lambda(y_\lambda^*)$.

Applying Corollary 4.3 then completes the proof. \square

PROPOSITION 8.2. *Assume that there exist functions $f(\cdot)$, $g(\cdot)$, H_λ , and J_λ such that, for any function $k_\lambda > 0$,*

$$F_\lambda(k_\lambda) \stackrel{\sim}{\sim} f(k_\lambda)H_\lambda \quad (29)$$

and

$$G_\lambda(k_\lambda) \stackrel{\sim}{\sim} g(k_\lambda)H_\lambda J_\lambda, \quad (30)$$

with $J_\lambda \gg 1$, and $g(\lambda^\alpha)J_\lambda Q(\lambda^\alpha) \ll f(\lambda^\alpha)$ for some $\alpha < 1/6$, so that certainly $F_\lambda(\kappa) \ll G_\lambda(\kappa)$ for all $\kappa > 0$.

Define

$$y_\lambda^* = \arg \min_{y>0} C[y; f, Q, gJ_\lambda]. \quad (31)$$

The staffing function y_λ^* is then asymptotically optimal in the sense of Corollary 4.3.

PROOF. In a similar fashion as in the proof of Theorem 8.1, it may be shown that $\lim_{\lambda \rightarrow \infty} x_\lambda^* = \infty$ and $\lim_{\lambda \rightarrow \infty} y_\lambda^* = \infty$. However, to show that the function $Q(\cdot)$ is an asymptotically exact approximation for the function $\pi_\lambda(\cdot)$ (using Lemma 5.2), we need to prove, additionally, that x_λ^* and y_λ^* do not grow faster than λ^α .

We start with showing that $x_\lambda^* \stackrel{\sup}{\leq} \lambda^\alpha$. Suppose to the contrary. Then, $x_\lambda^* \stackrel{\sup}{>} \lambda^\alpha$, so that $F_\lambda(x_\lambda^*) \stackrel{\sup}{>} F_\lambda(\lambda^\alpha)$ from (14). Using Lemma 5.2 and (29), (30),

$$\begin{aligned} F_\lambda(\lambda^\alpha) &\stackrel{\sim}{\sim} F_\lambda(\lambda^\alpha) + Q(\lambda^\alpha)G_\lambda(\lambda^\alpha) \\ &\stackrel{\sim}{\sim} F_\lambda(\lambda^\alpha) + \pi_\lambda(\lambda^\alpha)G_\lambda(\lambda^\alpha) = C_\lambda(\lambda^\alpha). \end{aligned}$$

Combining the above relations, we obtain $C_\lambda(x_\lambda^*) \geq F_\lambda(x_\lambda^*) \stackrel{\sup}{>} C_\lambda(\lambda^\alpha)$, contradicting the optimality of x_λ^* . Thus, $x_\lambda^* \stackrel{\sup}{\leq} \lambda^\alpha$. By similar arguments, $y_\lambda^* \stackrel{\sup}{\leq} \lambda^\alpha$. Hence, according to Lemma 5.2, $\pi_\lambda(x_\lambda^*) \stackrel{\sim}{\sim} \hat{\pi}_\lambda(x_\lambda^*) = Q(x_\lambda^*)$, and $\pi_\lambda(y_\lambda^*) \stackrel{\sim}{\sim} \hat{\pi}_\lambda(y_\lambda^*) = Q(y_\lambda^*)$. Applying Corollary 4.3 then completes the proof. Note that $F_\lambda(x_\lambda^*) \stackrel{\sim}{\sim} \hat{F}_\lambda(x_\lambda^*) = f(x_\lambda^*)H_\lambda$, $F_\lambda(y_\lambda^*) \stackrel{\sim}{\sim} \hat{F}_\lambda(y_\lambda^*) = f(y_\lambda^*)H_\lambda$, $G_\lambda(x_\lambda^*) \stackrel{\sim}{\sim} \hat{G}_\lambda(x_\lambda^*) = g(x_\lambda^*)H_\lambda J_\lambda$, and $G_\lambda(y_\lambda^*) \stackrel{\sim}{\sim} \hat{G}_\lambda(y_\lambda^*) = g(y_\lambda^*)H_\lambda J_\lambda$. \square

REMARK 8.3. The rationalized staffing function (20) is, in fact, also asymptotically optimal in the quality-driven regime as defined by (28) when $x_\lambda^* \ll \lambda^{1/6}$, since the proof of Proposition 8.2 then shows that $\pi_\lambda(x_\lambda^*) \stackrel{\sim}{\sim} Q(x_\lambda^*)$, while $Q(x_\lambda^*) \stackrel{\sim}{\sim} P(x_\lambda^*)$. However, the staffing function (31) is asymptotically optimal as well, while simpler.

EXAMPLE 8.4. Consider the same cost structure as in Example 6.3. Now suppose, however, that $a_\lambda \stackrel{\sim}{\sim} aJ_\lambda$, $b_\lambda \stackrel{\sim}{\sim} bJ_\lambda/\sqrt{\lambda}$, $c_\lambda \stackrel{\sim}{\sim} c$, and $d_\lambda \stackrel{\sim}{\sim} d/\sqrt{\lambda}$, with $J_\lambda \gg 1$. Then (29)–(30) are satisfied for $f(\kappa) = c\kappa$, $g(\kappa) = (a/\kappa) + b\sqrt{\mu}e^{-d\kappa\sqrt{\mu}}$, and $H_\lambda = \sqrt{\lambda/\mu}$. The asymptotically optimal staffing function is

$$y_\lambda^* = \arg \min_{y>0} \left\{ cy + Q(y) \left[\frac{a}{y} + b\sqrt{\mu}e^{-dy\sqrt{\mu}} \right] J_\lambda \right\}.$$

In the special case where the nonlinear penalty cost is asymptotically negligible, i.e., “ $b = 0$ ” (or rather, $b_\lambda\sqrt{\lambda}e^{-\mu\kappa d_\lambda\sqrt{\lambda/\mu}} \ll J_\lambda$), the optimal staffing functions reduces to

$$y_\lambda^* = \arg \min_{y>0} \left\{ cy + \frac{aQ(y)}{y} J_\lambda \right\}.$$

In the extreme case where the linear waiting costs are asymptotically insignificant, i.e., “ $a = 0$ ” (formally, $a_\lambda \ll J_\lambda$), the optimal staffing function takes the form

$$y_\lambda^* = \arg \min_{y>0} \{ cy + bQ(y)\sqrt{\mu}e^{-dy\sqrt{\mu}} J_\lambda \}.$$

9. Constraint Satisfaction

In the previous sections, we considered the problem of determining the staffing level to minimize the total staffing and waiting cost. A closely related problem, which is in fact motivated by actual practice, is to minimize the staffing level subject to a constraint $M_\lambda > 0$ on the waiting cost.

We are now interested in determining

$$N_\lambda^* := \min_{N > \lambda/\mu} \{N : K(N, \lambda) \leq M_\lambda\}, \quad (32)$$

with $K(N, \lambda) := \lambda\pi(N, \lambda/\mu)G(N, \lambda)$ denoting the waiting cost. Notice that $\lim_{N \rightarrow \infty} K(N, \lambda) = 0$, so N_λ^* is well defined.

As for the cost minimization problem, our approach is first to translate the discrete problem (32) into a continuous one, and then approximate the latter problem by a related continuous problem, which is easier to solve. Denote $x_\lambda^* := \min_{x > 0} \{x : K_\lambda(x) \leq M_\lambda\}$, with $K_\lambda(x) := \pi_\lambda(x)G_\lambda(x)$. Since $\pi_\lambda(\cdot)$ and $G_\lambda(\cdot)$ are both continuous and strictly decreasing, x_λ^* is the unique solution to the equation $K_\lambda(x) = M_\lambda$. Further, notice that $N_\lambda^* = \lceil N_\lambda(x_\lambda^*) \rceil$, which establishes the link between the discrete problem and the corresponding continuous problem.

To approximate x_λ^* , define z_λ^* as the solution to the equation $\hat{\pi}_\lambda(z)\hat{G}_\lambda(z) = M_\lambda$. The functions $\hat{\pi}_\lambda(\cdot)$ and $\hat{G}_\lambda(\cdot)$ that we consider will always be such that z_λ^* exists and is unique. We now formulate the approximation principle underlying our approach, in parallel to that for the cost minimization problem.

Define

$$T_\lambda(x) := \min\{|K(\lfloor N_\lambda(x) \rfloor, \lambda) - M_\lambda|, |K(\lceil N_\lambda(x) \rceil, \lambda) - M_\lambda|, |K(\lceil N_\lambda(x) \rceil, \lambda) - K(N_\lambda^*, \lambda)|\}. \quad (33)$$

LEMMA 9.1 (ASYMPTOTIC OPTIMALITY). *Denote $\hat{K}_\lambda(y) = \hat{\pi}_\lambda(y)\hat{G}_\lambda(y)$. Let z_λ^* be as defined above. If $K_\lambda(z_\lambda^*) \sim \hat{K}_\lambda(z_\lambda^*)$, then the staffing function z_λ^* is asymptotically optimal in the sense that, as $\lambda \rightarrow \infty$, $T_\lambda(z_\lambda^*) \ll M_\lambda$, with $T_\lambda(\cdot)$ given in (33).*

PROOF. For fixed λ , we distinguish between three cases.

(i) $N_\lambda^* - 1 < N_\lambda(z_\lambda^*) \leq N_\lambda^*$. Then $\lceil N_\lambda(z_\lambda^*) \rceil = N_\lambda^*$, so that $T_\lambda(z_\lambda^*) = 0$.

(ii) $N_\lambda(z_\lambda^*) \leq N_\lambda^* - 1$. Then $M_\lambda = K_\lambda(x_\lambda^*) \leq K(N_\lambda^* - 1, \lambda) \leq K(\lceil N_\lambda(z_\lambda^*) \rceil, \lambda) \leq K_\lambda(z_\lambda^*)$, so that $T_\lambda(z_\lambda^*) \leq |K(\lceil N_\lambda(z_\lambda^*) \rceil, \lambda) - M_\lambda| \leq K_\lambda(z_\lambda^*) - M_\lambda$.

(iii) $N_\lambda(z_\lambda^*) > N_\lambda^*$. Then $K_\lambda(z_\lambda^*) \leq K(\lfloor N_\lambda(z_\lambda^*) \rfloor, \lambda) \leq K(N_\lambda^*, \lambda) \leq K_\lambda(x_\lambda^*) = M_\lambda$, so that $T_\lambda(z_\lambda^*) \leq |K(\lfloor N_\lambda(z_\lambda^*) \rfloor, \lambda) - M_\lambda| \leq M_\lambda - K_\lambda(z_\lambda^*)$.

Thus, for all λ , $T_\lambda(z_\lambda^*) \leq |K_\lambda(z_\lambda^*) - M_\lambda| \ll M_\lambda$, as $\hat{K}_\lambda(z_\lambda^*) = M_\lambda$ by definition. \square

In full generality, it seems difficult to establish a stronger optimality property than indicated in the above lemma. Under additional conditions, however, it is possible to make sharper statements. For example, a more

desirable criterion for asymptotic optimality would be $|K(\lceil N_\lambda(z_\lambda^*) \rceil, \lambda) - K(N_\lambda^*, \lambda)| \ll M_\lambda$. And indeed, following the same reasoning as the proof of Lemma 9.1, it can be guaranteed to hold but under additional constraints on the oscillation of our costs.

9.1. Rationalized Regime

We first consider a rationalized scenario, by which we mean that, for some $\kappa > 0$,

$$G_\lambda(\kappa) \approx M_\lambda, \quad (34)$$

($\limsup_{\lambda \rightarrow \infty} (G_\lambda(\kappa)/M_\lambda) < \infty$ is actually sufficient) or in words, the waiting cost $G_\lambda(\cdot)$ is comparable to the constraint M_λ , as $\lambda \rightarrow \infty$.

For any $\lambda > 0$, define y_λ^* as the solution to the equation $P(y)G_\lambda(y) = M_\lambda$. Note that $P(\cdot)$ and $G_\lambda(\cdot)$ are both continuous and strictly decreasing with $\lim_{x \downarrow 0} P(x)G_\lambda(x) = \infty$ and $\lim_{x \rightarrow \infty} P(x)G_\lambda(x) = 0$, so that y_λ^* exists and is unique.

THEOREM 9.2. *The staffing function y_λ^* is asymptotically optimal in the sense of Lemma 9.1.*

PROOF. The idea of the proof may be described as follows. In order for Lemma 9.1 to apply, we need to show that the function $P(\cdot)$ is asymptotically close to the function $\pi_\lambda(\cdot)$ around y_λ^* . In view of Lemma 5.1, it suffices to show that $\limsup_{\lambda \rightarrow \infty} y_\lambda^* < \infty$. We prove this by contradiction, arguing that if this were not the case, then the incurred waiting cost would be strictly smaller than the constraint value M_λ .

We start with showing that $\limsup_{\lambda \rightarrow \infty} y_\lambda^* < \infty$. Suppose to the contrary. Then $y_\lambda^* \gg \kappa$, so that $P(y_\lambda^*)G_\lambda(y_\lambda^*) \ll P(\kappa)G_\lambda(\kappa)$ from (18) and the fact that $G_\lambda(\cdot)$ is decreasing. Using (34), $P(\kappa)G_\lambda(\kappa) \leq G_\lambda(\kappa) \approx M_\lambda$.

Combining the above relations, we deduce $P(y_\lambda^*) \times G_\lambda(y_\lambda^*) \ll M_\lambda$, contradicting the definition of y_λ^* . Thus, $\limsup_{\lambda \rightarrow \infty} y_\lambda^* < \infty$. Hence, according to Lemma 5.1, $\pi_\lambda(y_\lambda^*) \approx \hat{\pi}_\lambda(y_\lambda^*) = P(y_\lambda^*)$. Applying Lemma 9.1 then completes the proof. \square

PROPOSITION 9.3. *Assume that there exists a function $g(\cdot)$, such that, for any function $k_\lambda > 0$,*

$$G_\lambda(\kappa) \approx g(\kappa)M_\lambda, \quad (35)$$

so certainly $G_\lambda(\kappa) \approx M_\lambda$ for all $\kappa > 0$. Define y^ as the solution to the equation $P(y)g(y) = 1$. Assume that $g(\cdot)$ is continuous and decreasing, with $\lim_{x \downarrow 0} g(x) > 1$ so that y^* exists and is unique. The staffing function y^* is then asymptotically optimal in the sense of Lemma 9.1.*

PROOF. Similar to that of Theorem 9.2. Note that $G_\lambda(y^*) \approx \hat{G}_\lambda(y^*) = g(y^*)M_\lambda$. \square

EXAMPLE 9.4. Assume there is a waiting cost a_λ per customer per time unit, as well as a fixed penalty cost b_λ when the waiting time exceeds d_λ time units, i.e., $D_\lambda(t) = a_\lambda t + b_\lambda \mathbf{1}_{\{t > d_\lambda\}}$, so that $G(N, \lambda) = a_\lambda/(N\mu - \lambda) + b_\lambda e^{-(N\mu - \lambda)d_\lambda}$. Thus $G_\lambda(\kappa) = a_\lambda \sqrt{\lambda/\mu/\kappa} + b_\lambda \lambda e^{-\mu\kappa d_\lambda \sqrt{\lambda/\mu}}$.

(i) First, suppose that $a_\lambda \approx a\sqrt{\lambda/\mu}$, $b_\lambda e^{-\mu\kappa d_\lambda\sqrt{\lambda/\mu}} \ll 1$, and $M_\lambda = M\lambda$. Then, (35) is satisfied for $g(\kappa) = a/\kappa\mu M$. Proposition 9.3 says that the staffing function y^* determined as the unique solution to the equation $aP(y)/\mu y = M$ is asymptotically optimal in the sense of Lemma 9.1.

(ii) Now, suppose that $a_\lambda/\sqrt{\lambda} \ll 1$, $b_\lambda \approx b$, $d_\lambda \approx d/\sqrt{\lambda}$, and $M_\lambda = M\lambda$ with $b > M$. Then, (35) is satisfied for $g(\kappa) = b e^{-d\kappa\sqrt{\mu}}/M$. The asymptotically optimal staffing function is the unique solution to the equation $bP(y)e^{-dy\sqrt{\mu}} = M$.

(iii) Finally, consider the “combined” case where all costs show up in the limit, i.e., suppose that $a_\lambda \approx a\sqrt{\lambda/\mu}$, $b_\lambda \approx b$, $d_\lambda \approx d/\sqrt{\lambda}$, and $M_\lambda = M\lambda$. Then, (35) is satisfied for $g(\kappa) = a/\kappa\mu M + b e^{-d\kappa\sqrt{\mu}}/M$. The asymptotically optimal staffing function is the unique solution to the equation $P(y)[a/\mu y + b e^{-dy\sqrt{\mu}}] = M$.

EXAMPLE 9.5. An important special case is $a_\lambda = 0$, $b_\lambda = 1$, $d_\lambda = 0$, $M_\lambda = \epsilon\lambda$, which corresponds to a target waiting probability ϵ . Case (ii) of the above example then shows that the staffing function $y^* = P^{-1}(\epsilon)$ is asymptotically optimal. We described this example in §2—see (5). It is to be compared with (6), used in Whitt (1992) and Kolesar and Green (1998). On the differences and similarities between the two approximations, see §10.

9.2. Efficiency-Driven Regime

We now consider an efficiency-driven scenario, meaning that, for all $\kappa > 0$,

$$G_\lambda(\kappa) \ll M_\lambda, \quad (36)$$

(in fact $G_\lambda(\kappa) \leq \sup M_\lambda$ would be sufficient for the results below to hold) or in words, the waiting cost is dominated by the target upper bound, as $\lambda \rightarrow \infty$. For any $\lambda > 0$, define y_λ^* as the solution to the equation $G_\lambda(y) = M_\lambda$. Note that $G_\lambda(\cdot)$ is continuous and strictly decreasing with $\lim_{x \downarrow 0} G_\lambda(x) = \infty$ and $\lim_{x \rightarrow \infty} G_\lambda(x) = 0$, so that y_λ^* exists and is unique.

THEOREM 9.6. *The staffing function y_λ^* is asymptotically optimal in the sense of Lemma 9.1.*

PROOF. The idea of the proof is largely similar to that of Theorem 9.2. We start with showing that $\lim_{\lambda \rightarrow \infty} y_\lambda^* = 0$. Suppose to the contrary. Then $y_\lambda^* \sup > u$ for some $u > 0$, so that $P(y_\lambda^*)G_\lambda(y_\lambda^*) \inf < P(u)G_\lambda(u)$, from (17) and the fact that $G_\lambda(\cdot)$ is decreasing. Using (36), $P(u)G_\lambda(u) \leq G_\lambda(u) \sup \leq M_\lambda$.

Combining the above relations, we deduce $P(y_\lambda^*) \times G_\lambda(y_\lambda^*) \inf < M_\lambda$, contradicting the definition of y_λ^* . Thus, $\lim_{\lambda \rightarrow \infty} y_\lambda^* = 0$. Hence, according to Lemma 5.1, $\pi_\lambda(y_\lambda^*) \approx \hat{\pi}_\lambda(y_\lambda^*) = P(y_\lambda^*) \approx 1$. Applying Lemma 9.1 then completes the proof. \square

PROPOSITION 9.7. *Assume that there exist functions $g(\cdot)$ and J_λ such that, for any function $k_\lambda > 0$,*

$$G_\lambda(k_\lambda) \approx g(k_\lambda)J_\lambda M_\lambda, \quad (37)$$

with $J_\lambda \ll 1$, so certainly $G_\lambda(\kappa) \ll M_\lambda$ for all $\kappa > 0$. Define y_λ^* as the unique solution to the equation $g(y)J_\lambda = 1$. Assume that $g(\cdot)$ is continuous and strictly decreasing, with $\lim_{x \downarrow 0} g(x) = \infty$ so that y_λ^* exists and is unique. The staffing function y_λ^* is then asymptotically optimal in the sense of Lemma 9.1.

PROOF. Similar to that of Theorem 9.6. Note that $G_\lambda(y_\lambda^*) \approx g(y_\lambda^*)J_\lambda M_\lambda$. \square

EXAMPLE 9.8. Consider the same cost structure as in Example 9.4. Now suppose, however, that $a_\lambda \approx a\sqrt{\lambda/\mu}$, $b_\lambda \approx b$, $d_\lambda \approx d/\sqrt{\lambda}$, and $M_\lambda = M\lambda/J_\lambda$, with $J_\lambda \ll 1$. Then, (37) is satisfied for $g(\kappa) = a/\kappa\mu M + b e^{-d\kappa\sqrt{\mu}}/M$. The asymptotically optimal staffing function is the unique solution to the equation $[a/\mu y + b e^{-dy\sqrt{\mu}}]J_\lambda = M$.

In the special case where the nonlinear penalty cost is asymptotically negligible, i.e., “ $b = 0$ ” (or rather $b_\lambda e^{-\mu\kappa d_\lambda\sqrt{\lambda/\mu}} \ll 1$), the equation for the optimal staffing function reduces to $y_\lambda^* = (a/\mu M)J_\lambda$.

As noted earlier for the optimization problem, the extreme case where the linear waiting cost is asymptotically vanishing, i.e., “ $a = 0$ ” (formally, $a_\lambda\sqrt{\lambda} \ll 1$) does not make sense, since $\lim_{\kappa \downarrow 0} g(\kappa) = b/M < \infty$. Thus, the equation for the optimal staffing function will not have a positive solution for large λ .

9.3. Quality-Driven Regime

We finally consider a quality-driven scenario, meaning that, for all $\kappa > 0$,

$$G_\lambda(\kappa) \gg M_\lambda, \quad (38)$$

or in words, the waiting cost dominates the target upper bound, as $\lambda \rightarrow \infty$. For any $\lambda > 0$, define y_λ^* as the solution to the equation $G_\lambda(y)Q_\lambda(y) = M_\lambda$. Note that $G_\lambda(\cdot)$ and $Q_\lambda(\cdot)$ are continuous and strictly decreasing with $\lim_{x \downarrow 0} G_\lambda(x)Q_\lambda(x) = \infty$ and $\lim_{x \rightarrow \infty} G_\lambda(x)Q_\lambda(x) = 0$, so that y_λ^* exists and is unique.

THEOREM 9.9. *The staffing function y_λ^* is asymptotically optimal in the sense of Lemma 9.1.*

PROOF. The proof is largely similar to that of Theorem 9.2. However, now we rely on Lemma 5.2, which indicates that the function $Q_\lambda(\cdot)$ is asymptotically close to the function $\pi_\lambda(\cdot)$, provided $\lim_{\lambda \rightarrow \infty} y_\lambda^* = \infty$. This may be shown again by contradiction, verifying that if this were not the case, then the incurred waiting cost would now strictly exceed the constraint value M_λ .

We start with showing that $\lim_{\lambda \rightarrow \infty} y_\lambda^* = \infty$. Suppose to the contrary. Then $y_\lambda^* \inf < u$ for some $u > 0$, so that $P(y_\lambda^*)G_\lambda(y_\lambda^*) \sup > P(u)G_\lambda(u)$ from (16) and the fact that $G_\lambda(\cdot)$ is decreasing. Using Lemma 5.2 and (38), $P(u)G_\lambda(u) \approx P(u)G_\lambda(u) \gg M_\lambda$.

Table 2. Overview of the parameter settings for the numerical experiments.

$\mu = 1, c_\lambda = c = 1$						
Example 6.3						
i		ii		iii	Example 7.4	Example 8.4
$\lambda = 100$	$a = 2$	$\lambda = 100, d = 0.1$	$b = 5, d = 1$	$a = 2, b = 2.5, d = 0.1$	$a = 1, b_\lambda = 0$	$a = 1, b_\lambda = 0$
a	λ	b	λ	λ	λ	λ
0.1	5	0.1	5	5	10	10
0.25	6	0.25	6	6	20	20
0.5	7	0.5	7	7	30	30
1	⋮	1	⋮	⋮	⋮	⋮
2	⋮	2	⋮	⋮	⋮	⋮
4	99	4	99	99	190	190
10	100	10	100	100	200	200

Combining the above relations, we deduce $P(y_\lambda^*)G_\lambda(y_\lambda^*) \stackrel{\text{sup}}{>} M_\lambda$, contradicting the definition of y_λ^* . Thus, $\lim_{\lambda \rightarrow \infty} y_\lambda^* = \infty$. Hence, according to Lemma 5.2, $\pi_\lambda(y_\lambda^*) \sim \hat{\pi}_\lambda(y_\lambda^*) = Q_\lambda(y_\lambda^*)$. Applying Lemma 9.1 then completes the proof. \square

PROPOSITION 9.10. Assume that there exist functions $g(\cdot)$ and J_λ such that, for any function $k_\lambda > 0$,

$$G_\lambda(k_\lambda) \stackrel{\infty}{\sim} g(k_\lambda)J_\lambda M_\lambda, \quad (39)$$

with $J_\lambda \gg 1$, and $g(\lambda^\alpha)J_\lambda Q(\lambda^\alpha) \ll 1$ for some $\alpha < 1/6$, so that certainly $G_\lambda(\kappa) \gg M_\lambda$ for all $\kappa > 0$. Define y_λ^* as the unique solution to the equation $Q(y)g(y)J_\lambda = 1$. Assume that $g(\cdot)$ is continuous and strictly decreasing, with $\lim_{x \downarrow 0} g(x) > 0$ so that y_λ^* exists and is unique. The staffing function y_λ^* is then asymptotically optimal in the sense of Lemma 9.1.

PROOF. It may easily be shown that $y_\lambda^* \stackrel{\text{sup}}{\leq} \lambda^\alpha$. Hence, according to Lemma 5.2, $\pi_\lambda(y_\lambda^*) \sim \hat{\pi}_\lambda(y_\lambda^*) = Q(y_\lambda^*)$. The proof is further similar to that of Theorem 9.9. Note that $G_\lambda(y_\lambda^*) \stackrel{\infty}{\sim} \hat{G}_\lambda(y_\lambda^*) = g(y_\lambda^*)M_\lambda J_\lambda$. \square

EXAMPLE 9.11. Consider the same cost structure as in Example 9.4. Now suppose, however, that $a_\lambda \stackrel{\infty}{\sim} a\sqrt{\lambda/\mu}$, $b_\lambda \stackrel{\infty}{\sim} b$, $d_\lambda \stackrel{\infty}{\sim} d/\sqrt{\lambda}$, and $M_\lambda = M\lambda/J_\lambda$, with $J_\lambda \gg 1$. Then, (39) is satisfied for $g(\kappa) = a/\kappa\mu M + be^{-d\kappa\sqrt{\mu}}/M$. The asymptotically optimal staffing function is the unique solution of the equation $Q(y)[a/\mu y + be^{-dy\sqrt{\mu}}]J_\lambda = M$.

In the special case where the nonlinear penalty cost is asymptotically negligible, i.e., “ $b = 0$ ” (or rather $b_\lambda e^{-\mu\kappa d_\lambda \sqrt{\lambda/\mu}} \ll 1$), the equation for the optimal staffing function reduces to $(aQ(y)/\mu y)J_\lambda = M$.

In the extreme case where the linear waiting cost is asymptotically insignificant, i.e., “ $a = 0$ ” (formally, $a_\lambda \ll \sqrt{\lambda}$), the equation for the optimal staffing function takes the form $bQ(y)e^{-dy\sqrt{\mu}}J_\lambda = M$.

REMARK 9.12. Notice that the results for the constraint satisfaction problem closely mirror those for the cost minimization problem. In fact, the two problems may be formally related as follows. Consider a strictly decreasing function $M(\cdot)$ on $(0, \infty)$ with $\lim_{x \downarrow 0} M(x) = \infty$ and $\lim_{x \rightarrow \infty} M(x) = 0$. Then the (unique) solution to the equation $M(x) = 1$ is $\arg \min_{x > 0} \{M(x) + 1/M(x)\}$. Thus, the solution to the constraint satisfaction problem $\pi_\lambda(x)G_\lambda(x) = M_\lambda$ may also be represented as $\arg \min_{x > 0} \{M_\lambda^2/\pi_\lambda(x)G_\lambda(x) + \pi_\lambda(x)G_\lambda(x)\}$, which has the form of the cost minimization problem. The relation thus established is only formal. We could not utilize it to derive the results for constraint satisfaction from the optimization results.

10. Numerical Experiments

In this section, we present the results of some numerical experiments that we carried out. The main purpose of the numerical experiments was to test the accuracy of the approximations that arise from our asymptotically optimal staffing levels. The numerical results indicate that the rationalized approximation performs exceptionally well in all regimes. By Remark 7.3 we know that the rationalized approximation is, in fact, asymptotically optimal in all regimes. On the other hand, the accuracy displayed by the rationalized approximation is astonishingly better than our rigorous results lead us to believe. Our first two experiments address Grassmann (1988) and Kolesar and Green (1998), which correspond to Examples 6.3 and 9.5, respectively.

Grassmann (1988, Table 3) calculates the optimal staffing level N^* for the M/M/N queue, with offered loads $R = 1, 3, 10, 30, 100$ and costs $r = a/c = 10, 20, 100, 200$. While the latter are rather extreme values, our approximation (1) is nevertheless accurate: it is exact in 7 cases and off by only 1 agent in the other 13 cases.

Kolesar and Green (1998, Table 1) use (6) to calculate N^* that achieves $\Pr\{\text{Wait} > 0\} = \epsilon$, for $\epsilon = 0.2, 0.1, 0.05, 0.025, 0.01, 0.001$, and offered loads $R = 2^m$, $m =$

Table 3. Overview of the “wrong-regime tests.”

$\mu = 1, a = 1, b_\lambda = 0, c_\lambda = c = 1$						
Scaling	a_λ	Approximation Regime	$n^* - \lambda/\mu$	# Correct	# Off by 1	# Off by 2
Efficiency	$a\lambda^{-1/2}$	Rationalized	$y^*(a\lambda^{-1/2})\sqrt{\lambda/\mu}$	19	1	0
Efficiency	$a\lambda^{-1/2}$	Quality	$y_\lambda^*(a/\lambda)\sqrt{\lambda/\mu}$	2	14	4
Quality	$a\sqrt{\lambda}$	Rationalized	$y^*(a\sqrt{\lambda})\sqrt{\lambda/\mu}$	14	6	0
Quality	$a\sqrt{\lambda}$	Efficiency	$\sqrt{a/\mu}\lambda^{3/4}$	Poor—overstaffing by 10%–15%		
Rationalized	a	Quality	$y_\lambda^*(a/\sqrt{\lambda})\sqrt{\lambda/\mu}$	8	12	0
Rationalized	a	Efficiency	$\sqrt{a/\mu}\sqrt{\lambda}$	3	8	9

0, 1, . . . , 10. The approximation (1) is superior to (6). This is to be expected in view of the theory that supports the former, while the latter is heuristically based. Indeed, for $\epsilon = 0.2$, (1) is exact for 9 cases and misses by 1 agent for the other 2. In contrast, (6) misses by up to 7 agents. But more significantly, the misses in staffing levels lead to misses in the target delay probabilities, off by 25%–75% in 8 out of the 11 cases.

The approximation (6) improves as ϵ decreases, until eventually it coincides with (1). This is understood as follows: small values of ϵ give rise to large y^* (quality-driven), for which $\bar{\Phi}(y) \approx \phi(y)/y \approx P(y)$.

We now turn to numerical experiments related to Examples 6.3, 7.4, and 8.4. In all cases we compared an approximation to the optimal staffing level obtained from the asymptotics with the exact optimal staffing level, which we obtained through a simple search procedure. (The unimodality of $C(N, \lambda)$ in N makes such a search simple.)

In all the examples, we use $c_\lambda = c = 1$ and $\mu = 1$. Once we set $c_\lambda = c$, taking $c = 1$ is without loss of generality because we can take this as the definition of the monetary unit. Similarly, taking $\mu = 1$ is also without loss of generality because we can take $1/\mu$ as the definition of the time unit. The parameter settings for the experiments are summarized in Table 2.

We first describe the numerical results related to Example 6.3. We considered three cases:

- (i) $a_\lambda = a, \quad b_\lambda = 0$
- (ii) $a_\lambda = 0, \quad b_\lambda = b/\sqrt{\lambda}, \quad d_\lambda = d/\sqrt{\lambda}$
- (iii) $a_\lambda = a, \quad b_\lambda = b/\sqrt{\lambda}, \quad d_\lambda = d/\sqrt{\lambda}$.

These correspond to the three cases in Example 6.3. First consider case (i). For $r > 0$, let

$$y^*(r) = \arg \min_{y>0} \left\{ y + \frac{rP(y)}{y} \right\}. \quad (40)$$

We plot $y^*(r)$, $0 \leq r \leq 10$, in Figure 1.

Let $n^* = \lambda/\mu + y^*(r)\sqrt{\lambda/\mu}$. The rationalized approximation for case (i) of Example 6.3 is obtained by rounding n^* to the nearest integer. (The asymptotic analysis gives no guidance on how to go from n^* to an integer staffing level. Preliminary numerical calculations comparing rounding up, rounding down, and rounding off showed that rounding off is generally superior. So all of our numerical results involve

rounding off. Of course, if rounding off n^* yields a value smaller than or equal to λ/μ , the staffing level must be increased to be strictly greater λ/μ to avoid an unstable system.)

To check this approximation we first tried $\lambda = 100$ and the seven different values of a indicated in Table 2. In all these cases rounding off n^* gave the exact optimal staffing level. We next set $a = 2$ and tried all integer values of λ between 5 and 100. Here, rounding off n^* is never off by more than 1 agent, and is usually exact (83 out of 96 cases).

For case (ii) we let $n^* = \lambda/\mu + y^*(b, d)\sqrt{\lambda/\mu}$, where $y^*(b, d) = \arg \min_{y>0} \{y + bP(y)e^{-dy}\}$.

Here, we first tried $\lambda = 100$ and $d = 0.1$, with the seven different values of b indicated in Table 2. Again, we found that in all these cases rounding off n^* gave the exact optimal staffing level. Next, we set $b = 5$ and $d = 1$ and tried all integer values between 5 and 100 for λ . Rounding off n^* is almost always exact (84 out of 96 times), and is never off by more than 1.

For case (iii) we let $n^* = \lambda/\mu + y^*(a, b, d)\sqrt{\lambda/\mu}$, where $y^*(a, b, d) = \arg \min_{y>0} \{y + P(y)[a/y + be^{-dy}]\}$.

Here, we set $a = 2, b = 2.5$, and $d = 0.1$, and tried all integer values between 5 and 100 for λ (see Table 2). Again, rounding off n^* is almost always exact (80 out of 96 cases), and is never off by more than 1.

For Example 7.4 we restricted our attention to $b_\lambda = 0$. We initially set

$$a_\lambda = a\lambda^{-1/2}. \quad (41)$$

Defining $y_\lambda^*(a) = \arg \min_{y>0} \{y + a/(y\sqrt{\lambda})\}$, we can solve explicitly to obtain $y_\lambda^*(a) = \sqrt{a}\lambda^{-1/4}$. We thus let $n^* = \lambda/\mu + \sqrt{a/\mu}\lambda^{1/4}$.

For the numerical test of this approximation we set $a = 1$ and tried all integer multiples of 10 between 10 and 200 for λ (see Table 2), using (41) to determine a_λ . Rounding off n^* to the nearest integer is almost always exact (19 out of 20 times), and is never off by more than 1.

In Example 8.4 we again restricted our attention to $b_\lambda = 0$. We took

$$a_\lambda = a\sqrt{\lambda}. \quad (42)$$

Let $n^* = \lambda/\mu + y_\lambda^*(a)\sqrt{\lambda/\mu}$, where

$$y_\lambda^*(a) = \arg \min_{y>0} \left\{ y + \frac{aQ(y)}{\sqrt{\lambda}} y \right\}, \quad (43)$$

and $Q(y)$ is given by (13).

For the numerical test of this approximation, we again set $a = 1$ and tried all integer multiples of 10 between 10 and 200 for λ (see Table 2), using (42) to determine a_λ . Once again, rounding off n^* to the nearest integer is almost always exact (16 out of the 20 times), and is never off by more than 1.

The above tests were run under favorable conditions: The asymptotic scaling of the parameters was known, and the approximation used was that associated with the regime corresponding to the parameter scaling. On the other hand, the results of the test are astoundingly good: Rounding n^* to the nearest integer is almost always exact, and is never off by more than 1. Note that these tests include values of λ that do not appear to be very large.

It is clear that more numerical testing is in order. There are two aims to this testing: (1) Find parameter values that “break” the approximations, and (2) determine if any of the asymptotic approximations is robust enough to work outside of its regime and/or determine rules of thumb for when each approximation should be used. (Indeed, this last point is central for obtaining a practically useful approximation.)

The additional testing takes the form of “wrong-regime” testing. In all these tests we set $c_\lambda = c = 1$, $\mu = 1$, $b_\lambda = 0$, and $a = 1$. We used all integer multiples of 10 between 10 and 200 for λ . The results of these tests are summarized in Table 3. The first “wrong-regime” test we conducted involved scaling the parameters as in the efficiency-driven regime with $a_\lambda = a\lambda^{-1/2}$, and using the approximation from the rationalized regime. Thus, our approximation rounded off $n^* = \lambda/\mu + y^*(a\lambda^{-1/2})\sqrt{\lambda/\mu}$ to obtain the staffing level. The approximation was exact in all but one case ($\lambda = 160$), where it was off by 1.

Using the same parameters as in the preceding example, we used the approximation from the quality-driven regime (as if $a_\lambda = a\sqrt{\lambda}$). Thus our approximation rounded off $n^* = \lambda/\mu + y_\lambda^*(a/\lambda)\sqrt{\lambda/\mu}$, where y_λ^* is given by (43), to obtain the staffing level. The approximation was good, but not as good as the rationalized regime: It was off by 1 in 14 cases and off by 2 in 4 cases. We next scaled the parameters as in the quality-driven regime with $a_\lambda = a\sqrt{\lambda}$, and used the approximation from the rationalized regime, rounding off $n^* = \lambda/\mu + y^*(a\sqrt{\lambda})\sqrt{\lambda/\mu}$, where $y^*(\cdot)$ is given by (40), to obtain the staffing level. The approximation was off by 1 in 6 cases and exact in the others.

Using the same parameters as in the preceding example, we used the approximation from the efficiency-driven regime (as if $a_\lambda = a\lambda^{-1/2}$). Thus our approximation rounded off $n^* = \lambda/\mu + \sqrt{a/\mu}\lambda^{3/4}$ to obtain the staffing level. This approximation did not perform well, leading to overstaffing of 10%–15%.

Next, we scaled the parameters as in the rationalized regime, with $a_\lambda = a$, using the approximation from the quality-driven regime (as if $a_\lambda = a\sqrt{\lambda}$). Thus our approximation rounded off $n^* = \lambda/\mu + y_\lambda^*(a/\sqrt{\lambda})\sqrt{\lambda/\mu}$, where $y_\lambda^*(\cdot)$ is given by (43), to obtain the staffing level. The

approximation here was exact in 8 cases and off by 1 in 12 cases.

Using the same parameters as in the preceding example, we used the approximation from the efficiency-driven regime (as if $a_\lambda = a\lambda^{-1/2}$). Thus our approximation rounded off $n^* = \lambda/\mu + \sqrt{a/\mu}\sqrt{\lambda}$ to obtain the staffing level. The approximation is not as good as that of the rationalized or quality-driven regime: 3 cases were exact, 8 were off by 1, and 9 were off by 2.

In the tests above that involve scaling parameters for the efficiency-driven and quality-driven regimes, we used $a_\lambda = a\lambda^{-1/2}$ and $a_\lambda = a\sqrt{\lambda}$, respectively. These regimes hold more generally with $a_\lambda = a\lambda^{-\alpha}$ and $a_\lambda = a\lambda^\alpha$, respectively, for $\alpha > 0$. In addition to $\alpha = 1/2$, we also tried values of $\alpha = 1/4$ and 1. The runs for $\alpha = 1$ involved multiples of 50 from 50 to 1,000 for λ . (The approximations used the same value of α as used to scale the actual parameters.) The results can be summarized as follows. The rationalized approximation was excellent: It was mostly exact, and when it was wrong it was typically off by 1 and never off by more than 3. The quality-driven approximation was good, but not as good as the rationalized approximation. (Even in the quality-driven regime!) The efficiency-driven solution was the worst of the three, and substantially overstaffed in the quality-driven regime. In the efficiency-driven regime with $\alpha = 1$, the efficiency-driven solution provided the infeasible solution of λ as the staffing level. Of course this can be corrected by simply requiring that the staffing level must be strictly greater than λ .

11. Future Research

There are a few directions of research that suggest themselves. First, we would like to explain theoretically the extreme accuracy of our approximations. This cannot be anticipated from the corresponding asymptotic approximations. Next, the call center environment enjoys many features that are not captured by the M/M/N (Erlang-C) model. Important examples are customer abandonment, time-varying arrival rates, nonexponential service times, and multiple skill classes. The goal is to incorporate such features into our framework, which we now elaborate on.

A justification of the square-root safety staffing principle has been pursued at three levels: heuristic, conceptual and explicit. The *heuristic* level, as in Kolesar and Green (1998), is based on an infinite-server approximation. The *conceptual* level, as in Halfin and Whitt (1981), is founded on a formal limit theorem, the limit taken as the number of servers increases with the offered load in a precise manner. These two levels motivate staffing levels of the form $N \approx R + y^*\sqrt{R}$. The *explicit* level, as in the present work, involves calculation of the constant y^* as a function of model parameters.

(i) For models with abandonment (Erlang-A perhaps would be an appropriate term), the conceptual level was analyzed in Garnett et al. (2002). Notably, abandonment

renders the queue always stable, hence y^* can also take negative values. The explicit level for Erlang-A is now under study by the present authors.

(ii) The heuristic approach for time-varying arrival rates proved successful in Jennings et al. (1996). The corresponding conceptual level can be pursued within the service network framework of Mandelbaum et al. (1998).

(iii) As for general service times, the M/G/N queue is not amenable to exact analysis, and letting $N \rightarrow \infty$ does not make things easier. Indeed, the accuracy of the standard multiserver approximations turns out questionable as N becomes large, which is the relevant regime for call centers. A key challenge is the calculation of the Halfin-Whitt delay function for the M/G/N queue. To this end, one could first attempt the M/PH/N queue, following Puhalskii and Reiman (2000).

(iv) A first study on staffing algorithms for multiskill scenarios may be found in Borst and Seri (1997).

In the present paper, we assumed that the offered traffic parameters are exactly known. We then focused on determining the amount of safety staffing needed to deal with stochastic variability only. In practice, however, the offered traffic forecasts are typically not completely accurate (Jongbloed and Koole 2001), and additional staffing may be required in view of the inherent uncertainty in the parameters. While the need for accurate forecasts becomes even more pronounced with the high-level of utilization advocated by the present work, we still feel that our analysis for known parameters constitutes an essential first step.

Finally one could perhaps use duality theory from mathematical programming to relate the optimization approach to constraint satisfaction. This could possibly add insight on the optimality criteria for constraint satisfaction, which should be sharpened.

Appendix A. Proof of Lemma 5.2

LEMMA 5.2 *For any function x_λ with $\lim_{\lambda \rightarrow \infty} x_\lambda = \infty$, $\pi_\lambda(x_\lambda) \stackrel{\sim}{\sim} Q_\lambda(x_\lambda)$.*

If also $x_\lambda \leq \lambda^{1/6}$, then $\pi_\lambda(x_\lambda) \stackrel{\sim}{\sim} Q(x_\lambda)$. If specifically $x_\lambda = \kappa\sqrt{\lambda/\mu}$ for some constant $\kappa > 0$, then

$$\pi_\lambda(x_\lambda) \stackrel{\sim}{\sim} \frac{1}{\kappa\sqrt{2\pi\lambda/\mu}(1+\kappa)} \left(\frac{e^\kappa}{(1+\kappa)^{1+\kappa}} \right)^{\lambda/\mu}.$$

PROOF. The first statement follows after some manipulations from the proof of Proposition 1 of Halfin and Whitt (1981). If $x_\lambda = \kappa\sqrt{\lambda/\mu}$ for some constant $\kappa > 0$, then $N_\lambda(x_\lambda) = (1+\kappa)\lambda/\mu$, $r_\lambda(x_\lambda) = 1/(1+\kappa)$, and $1-r_\lambda(x_\lambda) = \kappa/(1+\kappa)$, so that

$$Q_\lambda(x_\lambda) = \frac{1}{\kappa\sqrt{2\pi\lambda/\mu}(1+\kappa)} \left(\frac{e^\kappa}{(1+\kappa)^{1+\kappa}} \right)^{\lambda/\mu}.$$

We now prove the second statement. Using the Taylor series expansion

$$\begin{aligned} \log u &= \log(1 - (1 - u)) \\ &= - \sum_{m=1}^{\infty} \frac{(1 - u)^m}{m} = -(1 - u) - \sum_{m=2}^{\infty} \frac{(1 - u)^m}{m}, \end{aligned}$$

we obtain

$$\begin{aligned} &\frac{\exp\{N_\lambda(x)[1-r_\lambda(x)+\log r_\lambda(x)]\}}{\sqrt{2\pi N_\lambda(x)}(1-r_\lambda(x))} \\ &= \frac{\exp\{-N_\lambda(x)\sum_{m=2}^{\infty} m^{-1}(1-r_\lambda(x))^m\}}{\sqrt{2\pi N_\lambda(x)}(1-r_\lambda(x))} \\ &= \frac{\exp\{-N_\lambda(x)(1-r_\lambda(x))^2/2 - N_\lambda(x)\sum_{m=3}^{\infty} m^{-1}(1-r_\lambda(x))^m\}}{\sqrt{2\pi N_\lambda(x)}(1-r_\lambda(x))} \\ &= Q(x) \frac{x}{\sqrt{N_\lambda(x)}(1-r_\lambda(x))} \exp\{[x^2 - N_\lambda(x)(1-r_\lambda(x))^2]/2\} \\ &\quad \cdot \exp\left\{-N_\lambda(x) \sum_{m=3}^{\infty} \frac{(1-r_\lambda(x))^m}{m}\right\}. \end{aligned}$$

Thus it remains to be shown that

$$\begin{aligned} &\frac{x_\lambda}{\sqrt{N_\lambda(x_\lambda)}(1-r_\lambda(x_\lambda))} \exp\{[x_\lambda^2 - N_\lambda(x_\lambda)(1-r_\lambda(x_\lambda))^2]/2\} \\ &\quad \cdot \exp\left\{-N_\lambda(x_\lambda) \sum_{m=3}^{\infty} \frac{(1-r_\lambda(x_\lambda))^m}{m}\right\} \approx 1 \end{aligned}$$

if $x_\lambda \ll \lambda^{1/6}$.

Note that

$$\sqrt{N_\lambda(x_\lambda)}(1-r_\lambda(x_\lambda)) = \frac{x_\lambda\sqrt{\lambda/\mu}}{\sqrt{\lambda/\mu + x_\lambda\sqrt{\lambda/\mu}}} \approx x_\lambda,$$

and

$$N_\lambda(x_\lambda)(1-r_\lambda(x_\lambda))^2 = \frac{x_\lambda^2\lambda/\mu}{\lambda/\mu + x_\lambda\sqrt{\lambda/\mu}},$$

so that

$$\begin{aligned} 0 &\leq x_\lambda^2 - N_\lambda(x_\lambda)(1-r_\lambda(x_\lambda))^2 \\ &= \frac{x_\lambda^3\sqrt{\lambda/\mu}}{\lambda/\mu + x_\lambda\sqrt{\lambda/\mu}} \leq \frac{x_\lambda^3}{\sqrt{\lambda/\mu}} \ll 1. \end{aligned}$$

Finally,

$$\begin{aligned} 0 &\leq N_\lambda(x_\lambda) \sum_{m=3}^{\infty} \frac{(1-r_\lambda(x_\lambda))^m}{m} \leq N_\lambda(x_\lambda) \sum_{m=3}^{\infty} (1-r_\lambda(x_\lambda))^m \\ &= \frac{N_\lambda(x_\lambda)(1-r_\lambda(x_\lambda))^3}{r_\lambda(x_\lambda)} \leq \frac{x_\lambda^3}{\sqrt{\lambda/\mu}} \ll 1, \end{aligned}$$

which completes the proof. \square

Appendix B. Properties of $P(\cdot)$

LEMMA B.1. *The function $P(\cdot)$ is strictly convex decreasing.*

PROOF. The function $P(\cdot)$ may be written $P(x) = 1/(1 + U(x))$, with $U(x) := xe^{x^2/2}V(x)$, and $V(x) := \int_{-\infty}^x e^{-y^2/2} dy$.

Differentiating, $P'(x) = (-U'(x)/(1 + U(x))^2)$, and $P''(x) = (2U'(x)^2 - U''(x)(1 + U(x)))/(1 + U(x))^3$.

Observing that $V'(x) = e^{-x^2/2}$ and $V''(x) = -xe^{-x^2/2}$, $U'(x) = x + (x^2 + 1)e^{x^2/2}V(x)$, and $U''(x) = x^2 + 2 + (x^3 + 3x)e^{x^2/2}V(x)$.

Thus, $P(\cdot)$ is decreasing since $U'(x) > 0$ for any $x > 0$. Also, $P(\cdot)$ is strictly convex because for any $x > 0$,

$$\begin{aligned} & 2U'(x)^2 - U''(x)(1 + U(x)) \\ &= 2[x + (x^2 + 1)e^{x^2/2}V(x)]^2 \\ &\quad - [x^2 + 2 + (x^3 + 3x)e^{x^2/2}V(x)][1 + xe^{x^2/2}V(x)] \\ &= x^2 - 2 + [2x^3 - x]e^{x^2/2}V(x) \\ &\quad + [x^4 + x^2 + 2][e^{x^2/2}V(x)]^2 \\ &> x^2 - 2 + [x^4 + 2x^3 + x^2 - x + 2e^{x^2/2}V(x)]e^{x^2/2}V(x) \\ &> x^2 - 2 + [x^4 + 2x^3 + x^2 - x + 2\sqrt{\pi/2}]e^{x^2/2}V(x) \\ &> x^2 + [x^4 + 2x^3 + x^2 - x + 2(\sqrt{\pi/2} - 1)]e^{x^2/2}V(x) \\ &> [x^2 - x + 2(\sqrt{\pi/2} - 1)]e^{x^2/2}V(x) \\ &= [(x - 1/2)^2 + 2(\sqrt{\pi/2} - 9/8)]e^{x^2/2}V(x) > 0. \quad \square \end{aligned}$$

Appendix C. Properties of $G_\lambda(\cdot)$

LEMMA C.1. *The function $G_\lambda(\cdot)$ is strictly convex decreasing.*

PROOF. It suffices to show that $K_\lambda(\cdot)$ is strictly convex decreasing with $K_\lambda(\omega) := \omega \int_0^\infty D_\lambda(t)e^{-\omega t} dt$.

Differentiating, $K'_\lambda(\omega) = \int_0^\infty [1 - \omega t]e^{-\omega t} D_\lambda(t) dt$, and $K''_\lambda(\omega) = \int_0^\infty [\omega t - 2]te^{-\omega t} D_\lambda(t) dt$, for all $\omega > 0$.

Since $D_\lambda(\cdot)$ is strictly increasing, we have $[1 - \omega t] \times D_\lambda(t) \leq [1 - \omega t]D_\lambda(\omega)$ for all $\omega > 0$, $t > 0$, with strict inequality for $t \neq 1/\omega$, so that

$$\begin{aligned} & \int_0^\infty [1 - \omega t]e^{-\omega t} D_\lambda(t) dt \\ & < D_\lambda(1/\omega) \int_0^\infty [1 - \omega t]e^{-\omega t} dt = 0, \end{aligned}$$

and

$$\begin{aligned} & \int_0^\infty [\omega t - 2]te^{-\omega t} D_\lambda(t) dt \\ & > D_\lambda(2/\omega) \int_0^\infty [\omega t - 2]te^{-\omega t} dt = 0, \end{aligned}$$

for all $\omega > 0$. \square

Acknowledgments

The authors gratefully acknowledge useful comments from the anonymous referee and associate editor, which helped improve the presentation of the results. The research of Avi Mandelbaum was carried out at Bell Labs, Lucent Technologies, and the Technion. The hospitality of the Mathematical Sciences Research Center at Bell Labs is greatly appreciated. At the Technion, the research was supported by the Israeli Science Foundation (ISF) grants 388/99 and 126/02, by the Niderzaksen Fund, and by the Technion funds for the promotion of research and sponsored research.

References

- Andrews, B., H. Parsons. 1993. Establishing telephone-agent staffing levels through economic optimizations. *Interfaces* **23**(2) 14–20.
- Borst, S. C., P. F. Seri. 1997. Robust algorithms for staffing agents with multiple skills. Unpublished manuscript.
- De Bruijn, N. G. 1981. *Asymptotic Methods in Analysis*. Dover, New York.
- Call Center Statistics (website). 2000. Retrieved September 2000, <http://www.callcenternews.com/resources/statistics.shtml>.
- Cooper, R. B. 1981. *Introduction to Queueing Theory*, 2nd ed. North Holland, New York.
- Erlang, A. K. 1948. On the rational determination of the number of circuits. E. Brockmeyer, H. L. Halstrom, A. Jensen, eds. *The Life and Works of A. K. Erlang*. The Copenhagen Telephone Company, Copenhagen, Denmark.
- Garnett, O., A. Mandelbaum, M. I. Reiman. 2002. Designing a telephone call center with impatient customers. *Manufacturing Service Oper. Management* **4** 208–227.
- Grassmann, W. K. 1986. Is the fact that the emperor wears no clothes a subject worthy of publication? *Interfaces* **16**(2) 43–51.
- Grassmann, W. K. 1988. Finding the right number of servers in real-world queueing systems. *Interfaces* **18**(2) 94–104.
- Halfin, S., W. Whitt. 1981. Heavy-traffic limits for queues with many exponential servers. *Oper. Res.* **29** 567–588.
- Jagers, A. A., E. A. Van Doorn. 1986. On the continued Erlang loss function. *Oper. Res. Lett.* **5** 43–46.
- Jagers, A. A., E. A. Van Doorn. 1991. Convexity of functions which are generalizations of the Erlang loss function and the Erlang delay function. *SIAM Rev.* **33** 281–282.
- Jennings, O. B., A. Mandelbaum, W. A. Massey, W. Whitt. 1996. Service staffing to meet time-varying demand. *Management Sci.* **42** 1383–1394.
- Jongbloed, G., G. M. Koole. 2001. Managing uncertainty in call centres using Poisson mixtures. *Appl. Stoch. Model Bus. Indust.* **17** 307–318.
- Kolesar, P. J., L. V. Green. 1998. Insights on service system design from a normal approximation to Erlang's delay formula. *Prod. Oper. Management* **7** 282–293.
- Mandelbaum, A., W. A. Massey, M. I. Reiman. 1998. Strong approximations for Markovian service networks. *Queueing Systems* **30** 149–201.
- Mandelbaum, A., W. A. Massey, M. I. Reiman, B. Rider, A. Stolyar. 2002. Queue length and waiting times for multiserver queues with abandonment and retrials. *Telecommunications Systems* **21** 149–172.
- Puhalskii, A. A., M. I. Reiman. 2000. The multiclass GI/PH/N queue in the Halfin-Whitt regime. *Adv. Appl. Probab.* **32** 564–595.
- Whitt, W. 1992. Understanding the efficiency of multi-server service systems. *Management Sci.* **38** 708–723.