

Service-Level Differentiation in Many-Server Service Systems via Queue-Ratio Routing

Itai Gurvich

Kellogg School of Management, Northwestern University, Evanston, Illinois 60208,
i-gurvich@kellogg.northwestern.edu

Ward Whitt

Department of Industrial Engineering and Operations Research, Columbia University, New York, New York 10027, ww2040@columbia.edu

Motivated by telephone call centers, we study large-scale service systems with multiple customer classes and multiple agent pools, each with many agents. To minimize staffing costs subject to service-level constraints, where we delicately balance the service levels (SLs) of the different classes, we propose a family of routing rules called *fixed-queue-ratio* (FQR) rules. With FQR, a newly available agent next serves the customer from the head of the queue of the class (from among those he is eligible to serve) whose queue length most exceeds a specified proportion of the total queue length. The proportions can be set to achieve desired SL targets. The FQR rule achieves an important *state-space collapse* (SSC) as the total arrival rate increases, in which the individual queue lengths evolve as fixed proportions of the total queue length. In the current paper we consider a variety of service-level types and exploit SSC to construct asymptotically optimal solutions for the staffing-and-routing problem. The key assumption in the current paper is that the service rates depend only on the agent pool.

Subject classifications: queues, networks: multiple classes, server pools; queues, optimization: design, staffing, routing; queues, limit theorems: asymptotic optimality, many-server heavy-traffic limits.

Area of review: Stochastic Models.

History: Received January 2007; revisions received October 2007, September 2008, January 2009; accepted April 2009.

Published online in *Articles in Advance*.

1. Introduction

Large call centers usually serve multiple classes of customers having different service requirements and different perceived value. The services provided by the call center agents usually require special skills, but it is usually not possible or cost effective for all agents to have all skills. With current technology, call centers have the capability of routing calls to appropriate agents with the required skills, using some form of *skill-based routing* (SBR), but it remains challenging to perform SBR effectively; see §5 of Gans et al. (2003).

Call centers usually specify their operational objectives in the form of *quality-of-service* (QoS) *constraints*. Following common practice, we will focus on the *x-y service-level* (SL) *constraint*, which stipulates that $x\%$ of the calls should be answered within y seconds. We let the call center have different SL constraints for different customer classes; e.g., with both regular and VIP customers, we might aim to respond to 80% of regular customers within 30 seconds, but 80% of VIP customers within 10 seconds.

In this context, the total problem has three components: design, staffing, and routing. In the design phase, we start by grouping the customers into classes and the agents into service pools. (In doing so, we assume that the customers within classes are homogeneous, as are the agents within

pools.) Then we must decide what skills each pool should have, i.e., which classes they are allowed to serve. In the staffing phase, we must decide how many agents should be in each service pool. Finally, in the routing phase we must decide how the agents should be assigned to customers in real time.

The total problem is typically large and complex, so that it is unproductive to search for an optimal solution. Thus we look for a *good, simple solution*, that produces near-optimal performance in a relatively simple way. In particular, we hope to turn the large scale into an advantage instead of a disadvantage by finding relatively simple procedures that become more effective as the scale increases. Indeed, we want to find a relatively simple approach that is asymptotically optimal for specified problems as the scale increases. Our goal is to achieve *simplicity and asymptotic optimality*.

1.1. A Simple Intuitive Routing Rule: FQR

When considering possible controls, we think we should seek controls that are intuitive and structurally simple. Controls that lack any evident structure or insight are unlikely to be used by call center managers. A good example of a simple and intuitive control that is applicable to very general network structures (but essentially limited to

single-agent service pools) is the *generalized-c μ* ($Gc\mu$) rule, first introduced by Van Mieghem (1995) for a model with multiple classes and a single server (also known as the V model), and generalized to more complicated networks by Mandelbaum and Stolyar (2004). A parallel to Mandelbaum and Stolyar (2004) in a many-server setting has been provided by Atar (2005), who characterizes a family of controls that achieve asymptotically optimal performance in the QED regime. (See Gurvich and Whitt 2009b for more discussion.) Although the controls in Atar (2005) can be implemented easily in a computerized environment, they are not nearly as simple as the $Gc\mu$ rule. Thus, it seems desirable to seek a family of controls for many-server systems that bridge the gap between the simple and intuitive $Gc\mu$ rule in Mandelbaum and Stolyar (2004) and the more complicated controls in Atar (2005).

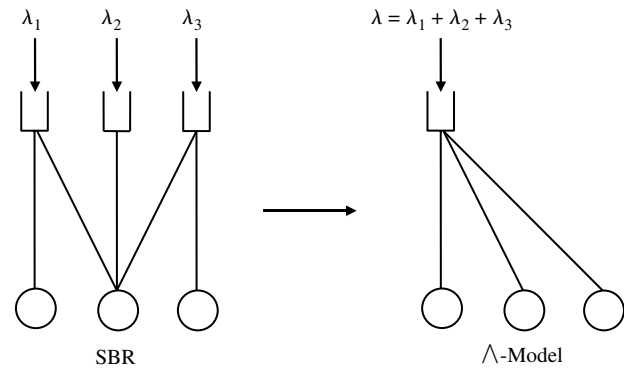
With that goal in mind, we propose *fixed-queue-ratio* (FQR) routing. We assume that there is a queue for each customer class. When an agent becomes free, he chooses the customer from the head of the line (from one of the classes he can serve) for which the queue length most exceeds a fixed proportion p_i of the total queue length (for all classes). The proportions p_i are in turn chosen to depend on the specified SL constraints. The FQR rule is a special case of the *queue-and-idleness-ratio* (QIR) family of controls that we introduce in Gurvich and Whitt (2009a). A consequence of Gurvich and Whitt (2009b) is that FQR makes the separate queue lengths *asymptotically proportional* to the total queue length. In other words, FQR produces a very important *state-space collapse* (SSC), causing the vector-valued queue-length process to evolve, asymptotically, as a one-dimensional process. In addition, FQR is a simple balancing rule like the $Gc\mu$ rule and, like the $Gc\mu$ rule, it is a highly decentralized control.

The key assumption that we make here is that the service rates are *pool dependent*, i.e., that the service time of a customer depends on the type (pool) of the agent that provides the service, but not on the customer class. Under this assumption, SSC reduces the multidimensional system dynamics to a tractable one-dimensional process. This reduction allows us to provide closed-form expressions for the staffing levels and prove that our proposed solution is not only feasible, but also asymptotically optimal.

Our solution stands on two pillars: (i) staffing via reduction of the multiclass multipool system to a single-class multipool system, and (ii) a simple routing rule that simultaneously makes the system perform as efficiently as the single-class multipool system and takes care of the service-level differentiation.

Our *reduction approach to staffing* illustrates our emphasis on simplicity: We propose first choosing the total number of agents by aggregating the service-level constraints and acting as if all customers have access to all agents. Thus, we reduce the original SBR system into a single-class multipool call center known as the inverted-V model; see Figure 1.

Figure 1. An SBR model and its corresponding \wedge model.



In the second step, we aim to satisfy the individual class-level QoS constraints by appropriately routing the customers in the system. Because the staffing and design decisions are highly interdependent, our proposed approach may seem naive and inadequate. However, we show that the FQR routing rule that we use in the second step guarantees that the two-step solution is nearly optimal, thus decomposing the joint optimization problem of design, staffing, and routing into more elementary problems that can be addressed sequentially.

In our sequel, Feldman et al. (2007), we remove the pool-dependence assumption and provide a general asymptotic feasibility (but not optimality) result. Based on the feasibility result, we then construct simple simulation-based optimization procedures to solve the design, staffing, and routing problem for more-general SBR systems.

1.2. Related Literature

Several recent papers have used the simplified reduction approach to staffing. Theoretical support is contained in Armony (2005) and Gurvich et al. (2008). These papers established asymptotic optimality of that staffing approach with appropriate routing for special classes of models as the total arrival rate increases. The first paper considered models with a single customer class and multiple agent types, whereas the second considered symmetric models with multiple customer classes, but a single agent pool. Their asymptotic optimality follows Borst et al. (2004), which formulated and established asymptotic optimality for the single-class, single-pool $M/M/N$ queue. The asymptotic framework is the now-familiar many-server heavy-traffic limiting regime introduced by Halfin and Whitt (1981), which is also known as the *quality-and-efficiency-driven* (QED) regime. In the QED regime the arrival rate and numbers of servers both increase, whereas the service-time distribution remains unchanged. These two limits are coordinated so that the probability of delay approaches a limit strictly between 0 and 1. Borst et al. (2004) showed that the QED regime arises naturally from economic considerations. We will be considering the QED regime throughout this paper.

The simplified reduction approach to staffing was also a central idea in Wallace and Whitt (2005), which developed a simulation-based iterative algorithm for staffing an SBR call center that starts by choosing an initial total number of agents by acting as if the call center were a single-class single-skill call center. After initial skill requirements are assigned, simulation is used iteratively to find detailed staffing and skill requirements so that the SL and other QoS constraints are met. The approach in Wallace and Whitt (2005) has two shortcomings, that we address here. First, that approach requires an iterative simulation algorithm to adjust staffing levels and skill assignments in order to satisfy the class-dependent QoS constraints. Because service is performed in a relatively short time scale compared to staffing, we think it should be more effective to primarily rely on the routing rather than the staffing in order to achieve desired service differentiation. In this paper we provide a way to do that. Second, although the approach in Wallace and Whitt (2005) seems to become more effective as the scale increases, it has not yet been shown to be asymptotically feasible or optimal as the scale increases. Here, in contrast, we establish asymptotic optimality. In this paper we assume that the service rates are *pool dependent*. A special case is the system with common service rates, e.g., as considered in Wallace and Whitt (2005).

The analysis in this paper relies heavily on our previous paper (Gurvich and Whitt 2009a), which establishes SSC results for the QIR generalization of FQR. We establish asymptotic optimality for the case of convex holding costs in Gurvich and Whitt (2009b).

An important contribution here is simultaneously addressing the three problems of design, staffing, and routing. Conventionally, these are treated separately and hierarchically. Wallace and Whitt (2005) also addressed these three problems together, but the only previous work we are aware of that establishes asymptotic feasibility or optimality for all three problems is Bassamboo et al. (2006a, b), Bassamboo and Zeevi (2008). They establish asymptotic optimality for the problem of minimizing costs associated with waiting, abandonments, and customer rejections. In Bassamboo and Zeevi (2008) they also consider abandonment constraints, but not tail-probability SL constraints. Their analysis is interesting because it focuses on uncertainty in the arrival rates. As a consequence, they consider a different limiting regime, the efficiency-driven (ED) regime. Their general setting allowing for uncertainty in the arrival rates comes at the price of having to restrict the analysis to a cruder notion of asymptotic optimality than the one we use here. Our finer analysis, although more limited in its scope, allows us to identify key system characteristics and, in turn, to construct intuitive routing schemes. Moreover, it allows us to tackle directly waiting-time tail-probability SL constraints that are widely used in the industry. In addition, the routing scheme we propose in this paper is used in Gurvich et al. (2008) to construct a staffing and routing algorithm

for a call center with uncertain arrival rates operating under SL constraints.

Within the context of single-server stations, several papers have tackled the problem of SL constraint satisfaction. Notable is Van Mieghem (2003), which embeds the constraint-satisfaction problem into the convex holding cost setting of Van Mieghem (1995), rather than dealing with it directly.

Organization of the Paper. In §2 we introduce the model and initial problem formulation. The proposed design, staffing, and routing solution is introduced in §3 and the asymptotic optimality results are stated in §4. We introduce and solve additional problem formulations in §5. We state conclusions in §6. Some proofs and auxiliary results appear in the e-companion, which is available as part of online version that can be found at <http://or.journal.informs.org/>.

2. The SBR System and the Problem Formulation

We consider a system with a set $\mathcal{I} := \{1, \dots, I\}$ of customer classes and a set $\mathcal{J} := \{1, \dots, J\}$ of agent types. The number of agents of type j (which will be a decision variable) is denoted by N_j ; let $N := (N_1, \dots, N_J)$. Class- i customers arrive according to a Poisson process with rate λ_i and $\lambda := \sum_{i \in \mathcal{I}} \lambda_i$ is the aggregate arrival rate. If a type- j agent can serve a class- i customer, we let $\mu_{i,j}$ be the corresponding service rate. Alternatively, the mean handling time of a class- i customer by a type- j agent is $1/\mu_{i,j}$. A *key assumption in this paper is that the service rates are pool dependent*. That is, for each $j \in \mathcal{J}$, we have $\mu_{i,j} = \mu_j$ for all $i \in \mathcal{I}$ that can be served by pool j . Throughout, we will assume, without loss of generality, that the service rates are labeled in decreasing order:

$$\mu_1 \geq \dots \mu_2 \geq \dots \geq \mu_J.$$

At this point, we do not consider customer abandonment; see §5.2 for extensions to that case.

The possible-routing graph for this SBR system has a natural representation as a bipartite graph (see Figure 1) with vertices $V = \mathcal{I} \cup \mathcal{J}$; i.e., V is the union of the set of customer classes and the set of agent pools. Then, the only edges we consider connect customer classes to agent pools: $E = \{(i, j) \in \mathcal{I} \times \mathcal{J} : j \text{ can serve } i\}$. An edge (i, j) is present in the routing graph if class- i customers can be served by type- j agents. We let $I(j)$ be the set of customer classes that can be served by pool j , i.e., $I(j) := \{i \in \mathcal{I} : (i, j) \in E\}$, and, symmetrically, we define $J(i) := \{j \in \mathcal{J} : (i, j) \in E\}$ to be the set of pools that are qualified to serve class- i customers. In addition, we assume that agents of type- j incur a cost of c_j per unit of time. We also will allow the system to impose additional constraints on the staffing vector to reflect union contracts, hiring and

training constraints, or other managerial considerations; see §2.1 for the formal modeling of these constraints.

For the asymptotic analysis, we will construct a sequence of SBR systems indexed by the aggregate arrival rate λ . The service rates μ_j and the routing graph are held fixed. We will make the dependence on the index explicit by adding the superscript λ to all the relevant parameters and processes. We assume that the ratios $a_i := \lambda_i/\lambda$ remain constant for all λ . Also, we let the SL target T_i^λ scale with λ to put the system into the QED regime.

ASSUMPTION 2.1 (QED SCALING FOR SL TARGETS). *The SL targets T_i^λ , $i \in \mathcal{F}$ are scaled so that $T_i^\lambda = \bar{T}_i/\sqrt{\lambda}$ for some strictly positive constants \bar{T}_i , $i \in \mathcal{F}$.*

The proposed staffing and routing solution will be completely defined in terms of the original targets T_i^λ ; there will be no use of the constants \bar{T}_i . The scaling is only used for the proof of asymptotic optimality.

2.1. The Problem Formulation

To formulate our optimization problem, let $A_i^\lambda(t)$ be the number of class- i customers to arrive by time t . Let $\bar{W}^{\lambda,T}$ be the average waiting time of all customers that arrived up to time T ; let $F_i^{\lambda,T}(\cdot)$ be the empirical distribution of the waiting time of class- i customers up to time T ; and let $\bar{F}_i^{\lambda,T}(\cdot)$ be its complement; i.e.,

$$\bar{W}^{\lambda,T} := \frac{\sum_{i=1}^I \sum_{k=1}^{A_i^\lambda(T)} w_{i,k}^\lambda}{A^\lambda(T)} \quad \text{and} \quad (1)$$

$$\bar{F}_i^{\lambda,T}(y) := \frac{\sum_{k=1}^{A_i^\lambda(T)} \mathbf{1}\{w_{i,k}^\lambda > y\}}{A_i^\lambda(T)},$$

where $A^\lambda(T) := \sum_{i \in \mathcal{F}} A_i^\lambda(T)$, $w_{i,k}^\lambda$ is the realized waiting time of the k th class- i customer to arrive to the system after time 0, and $\mathbf{1}B$ is the indicator of the event B , which is equal to 1 if B occurs, and 0 otherwise. An initial formulation, representing common call-center goals, can then be stated as follows:

$$\begin{aligned} & \text{minimize} \quad \sum_{j \in \mathcal{F}} c_j N_j \\ & \text{subject to} \quad \limsup_{T \rightarrow \infty} \bar{F}_i^{\lambda,T}(T_i^\lambda) \leq \alpha, \quad 1 \leq i \leq I, \quad (2) \\ & \quad \quad \quad N \in \mathcal{A}^\lambda, \quad \pi \in \Pi, \end{aligned}$$

where \mathcal{A}^λ is a subset of \mathbb{Z}_+^J that is generated by linear constraints. Specifically, we allow constraints of the form $N \in \mathcal{A}^\lambda := \mathcal{A}_b^\lambda$, where $\mathcal{A}_b^\lambda = \{N \in \mathbb{Z}_+^J : A \cdot N \leq b^\lambda\}$, for some matrix $A \in \mathbb{R}^{d \times J}$, $d \in \mathbb{Z}_+$, and $b^\lambda = \lambda \hat{b}$ for some $\hat{b} \in \mathbb{R}^d$. We also let $\tilde{\mathcal{A}}^\lambda := \tilde{\mathcal{A}}_b^\lambda$ be the set obtained from \mathcal{A}^λ by relaxing the integrality assumptions. That is, $\tilde{\mathcal{A}}^\lambda = \{N \in \mathbb{R}_+^J : A \cdot N \leq b^\lambda\}$. If the set of possible staffing vectors is unconstrained, we set $\mathcal{A}^\lambda \equiv \mathbb{Z}_+^J$.

We observe that to consider optimality, (2) is not well formulated, because we have not yet sufficiently constrained

the policies. So far, the formulation permits giving some customers satisfactory performance at the expense of giving other customers (in the proportion $1 - \alpha$) arbitrarily poor performance. This problem is discussed extensively at the end of §2 of Gurvich et al. (2008), so we will be brief here. To illustrate the difficulties, note that we could elect not to serve class- i customers who have waited longer than T_i . Even if we required first-come first-served (FCFS) service within each class, we could satisfy all the constraints with relatively limited staffing by disallowing any waiting, i.e., by using a pure-loss model. Clearly, in a loss model all the customers that do enter the system do not experience any wait, and we may choose the number of agents so that the blocking probability is less than $1 - \alpha$. That is clearly an undesirable outcome because many customers are blocked and do not receive service at all. Even when requiring that all customers be served, highly undesirable policies are possible, such as the alternating-priority control discussed by Gurvich et al. (2008).

As a consequence, we modify (2) by adding an additional constraint; in particular, we initially consider the following *best-effort optimization problem*:

$$\begin{aligned} & \text{minimize} \quad \sum_{j \in \mathcal{F}} c_j N_j \\ & \text{subject to} \quad \limsup_{T \rightarrow \infty} \bar{W}^{\lambda,T} \leq T_i^\lambda, \quad (3) \\ & \quad \quad \quad \limsup_{T \rightarrow \infty} \bar{F}_i^{\lambda,T}(T_i^\lambda) \leq \alpha, \quad 1 \leq i \leq I-1, \\ & \quad \quad \quad N \in \mathcal{A}^\lambda, \quad \pi \in \Pi. \end{aligned}$$

We emphasize that it is not sufficient to add the global average-waiting-time constraint. It is also important to remove the individual SL constraint of class I . Otherwise, the problems with formulation (2) remain unresolved. Indeed, if the global average-waiting-time constraint is very loose, one can show that it is possible to construct alternating priority controls as illustrated in Gurvich et al. (2008) to construct policies under which, part of the time, each class experiences extremely low service levels.

The difficulties in formulation (2) can be avoided in other ways, e.g., considering only average-waiting-time constraints. Such a formulation and its corresponding solution are considered in §5. We first focus on the formulation (3).

We now define the set of admissible policies Π . To this end, we say that *all customers are served* if it is not allowed to block or overflow customers; i.e., we require that for all $t \geq 0$, $Q_i(t) = A_i(t) - D_i(t) - Z_i(t)$, where $D_i(t)$ and $Z_i(t)$ are, respectively, the number of class- i departures from the system up to time t and the number of class- i customers in service at time t .

We say that a routing policy is *nonanticipative* if a decision at any time is based on the history up to that time and not upon future events. We say that a routing policy is *nonpreemptive* if customers stay in service with the agent first assigned to them until their service is complete once an agent has been assigned.

DEFINITION 2.1 (ADMISSIBLE ROUTING POLICIES). We say that a routing policy π is *admissible* if: (1) it is nonanticipative, (2) it is nonpreemptive, and (3) all customers are served. Let Π be the set of all admissible routing policies.

We conclude this section with the definition of asymptotic feasibility.

DEFINITION 2.2 (ASYMPTOTIC FEASIBILITY). A sequence of staffing vectors and routing policies $\{(N^\lambda, \pi^\lambda)\}$ is asymptotically feasible for (3) if (a) $\pi^\lambda \in \Pi$, for all λ , and (b) for every $\epsilon > 0$, there exists $T^*(\epsilon)$ such that, for all $T \geq T^*(\epsilon)$,

$$\limsup_{\lambda \rightarrow \infty} P\left\{\frac{\bar{W}^{\lambda, T}}{T_i^\lambda} \geq 1 + \epsilon\right\} \leq \epsilon, \quad (4)$$

and

$$\limsup_{\lambda \rightarrow \infty} P\{\bar{F}_i^{\lambda, T}(T_i^\lambda) \geq \alpha + \epsilon\} \leq \epsilon, \quad 1 \leq i \leq I - 1. \quad (5)$$

Asymptotic feasibility holds for (2) instead of (3) if (4) and (5) above are replaced by

$$\limsup_{\lambda \rightarrow \infty} P\{\bar{F}_i^{\lambda, T}(T_i^\lambda) \geq \alpha + \epsilon\} \leq \epsilon, \quad i \in \mathcal{F}. \quad (6)$$

Given two positive real-valued functions f and g , we say that $f(x)$ is $o(g(x))$ (as $x \rightarrow \infty$) if $f(x)/g(x) \rightarrow 0$ as $x \rightarrow \infty$; we say that $f(x)$ is $O(g(x))$ if $f(x)/g(x)$ is bounded as $x \rightarrow \infty$. The definition of asymptotic optimality will be the same for all the formulations that we consider in this paper. Hence, we do not specify a specific formulation within the definition. For the rest of the paper, asymptotic optimality will always be in the sense of Definition 2.3 below, where asymptotic feasibility will depend on the context. The $o(\sqrt{\lambda})$ in the asymptotic optimality condition (7) below corresponds to asymptotic optimality in the diffusion scale, which is more refined than asymptotic optimality in the fluid scale, which would involve a larger bound on the error of $o(\lambda)$ as $\lambda \rightarrow \infty$. (The staffing levels will be $O(\lambda)$.)

DEFINITION 2.3 (ASYMPTOTIC OPTIMALITY). A sequence of staffing vectors and routing policies $\{(N^\lambda, \pi^\lambda)\}$ is asymptotically optimal if it is asymptotically feasible, and

$$[c \cdot N^\lambda - c \cdot \tilde{N}^\lambda]^+ = o(\sqrt{\lambda}) \quad \text{as } \lambda \rightarrow \infty, \quad (7)$$

for any other sequence $\{(\tilde{N}^\lambda, \pi^\lambda)\}$ of asymptotically feasible staffing vectors and routing policies.

We end this section with a brief discussion of the relation between our notions of asymptotic feasibility and optimality, and the more traditional steady-state feasibility and optimality.

REMARK 2.1 (STEADY-STATE CONSTRAINTS). Although actual call-center operations involve finite-horizon decisions and constraints, the traditional way of call-center modeling would be to write both (2) and (3) with steady-state

constraints instead of the finite-horizon ones. To obtain the steady-state formulation, one would replace the individual tail constraints with $P\{W_i^\lambda(\infty) > T_i^\lambda\} \leq \alpha$, where $W_i^\lambda(\infty)$ is the class- i steady-state waiting time. The global average delay constraint is replaced with the constraint $E[W^\lambda(\infty)] \leq T_i^\lambda$ where $W^\lambda(\infty)$ is the steady-state global waiting time, i.e., $W^\lambda(\infty)$ is equal in distribution to $\sum_{i \in \mathcal{F}} (\lambda_i/\lambda) W_i^\lambda(\infty)$. Even though these alternative formulations differ little from a practical perspective, they are mathematically different. They are asymptotically equivalent only if one can establish a certain limit-interchange result. Such limit-interchange arguments are elementary for some models, such as the inverted-V model, but it is a complex task for the general SBR setting. In this paper we restrict the attention to the long-run average formulation in (2) and (3). What we do parallels what is done with simulation. \square

2.2. A Lower Bound

Our solution will be based on a reduction of the SBR system to a more elementary model in which multiple agent types serve a single customer class, also known as the inverted-V (or \wedge) model. Given an SBR system, the associated \wedge model has the same set of agent-pools \mathcal{F} , the same staffing levels $\{N_j, j \in \mathcal{F}\}$, and the same service rates $\{\mu_j, j \in \mathcal{F}\}$. In addition, the arrival rate of the single customer class is λ —the sum of the arrival rates in the SBR system. An example of an SBR system and its corresponding \wedge model is given in Figure 1.

Clearly, the \wedge model is not as simple as the $M/M/N$ queue. However, when it is optimally operated, its asymptotic performance leads to simple expressions for staffing, as has been shown by Armony (2005). We will exploit the results in Armony (2005) here. In particular, we will exploit a result for the \wedge model, stating that

$$E[W^\lambda(\infty)] \leq T_i^\lambda \quad \text{only if } \sum_j \mu_j N_j^\lambda \geq \lambda + \beta^* \sqrt{\lambda} + o(\sqrt{\lambda}), \quad (8)$$

where β^* is the unique solution to

$$\frac{\mathbf{P}_{\mu_1}(\beta)}{\beta} = \sqrt{\lambda} T_i^\lambda =: \bar{T}_i, \quad \text{and} \quad (9)$$

$$\mathbf{P}_{\mu_1}(\beta) := \left[1 + \frac{(\beta/\sqrt{\mu_1})\Phi(\beta/\sqrt{\mu_1})}{\phi(\beta/\sqrt{\mu_1})} \right]^{-1},$$

with $\phi(\cdot)$ and $\Phi(\cdot)$ being, respectively, the standard normal pdf and cdf. Here, $\mathbf{P}_{\mu_1}(\beta)$ is the asymptotic delay probability in the \wedge model operated under the fastest-server-first (FSF) policy as introduced by Armony (2005). That is, assuming that $\sum_{j \in \mathcal{F}} \mu_j N_j^\lambda = \lambda + \beta \sqrt{\lambda} + o(\sqrt{\lambda})$ and that FSF is used (plus additional technical conditions), Armony (2005) shows that $P\{W^\lambda(\infty) > 0\} \rightarrow \mathbf{P}_{\mu_1}(\beta)$, with $W^\lambda(\infty)$ being the steady-state waiting time in the \wedge model.

We note that the necessary condition (8) is given in Armony (2005) for steady-state asymptotic feasibility,

whereas we focus here on the somewhat weaker notion that appears in Definition 2.2. However, we find that the same necessary condition holds if one considers the same \wedge model, but with long-run average constraints (as in (3)) and with our notion of asymptotic feasibility. This is proved in Lemma EC.2.2 of the online appendix which, in turn, is a key step in the proof of Theorem 2.1 below.

We now use this necessary condition to construct a lower bound for the staffing of the SBR model. In doing this construction we will use two facts:

(i) In contrast to the SBR system, customers in the \wedge model have access to all agent pools. It is intuitively clear, then, that if a given staffing vector is not sufficient for the given aggregate waiting-time target in the \wedge model, it will also not be sufficient in the less efficient SBR system. Consequently, to meet the global waiting-time constraint it is necessary that $\sum_{j \in \mathcal{J}} \mu_j N_j^\lambda \geq \lambda + \beta^* \sqrt{\lambda} + o(\sqrt{\lambda})$.

(ii) Because we have no abandonments in the system, the capacity should suffice to serve all customers (at least at a fluid scale). In particular, any feasible staffing vector must satisfy that

$$\sum_{j \in J(i)} \mu_j N_j y_{i,j} \geq \lambda_i, \quad i \in \mathcal{F}$$

for some vector y such that $\sum_{i \in I(j)} y_{i,j} \leq 1$ and $y_{i,j}$, $(i, j) \in E$ are positive.

Together, these two informal arguments suggest that an asymptotic lower bound for the optimization problem (3) should be given by

$$\begin{aligned} & \text{Minimize} \quad \sum_{j \in \mathcal{J}} c_j N_j \\ & \text{Subject to:} \quad \sum_{j \in \mathcal{J}} \mu_j N_j \geq \lambda + \beta^* \sqrt{\lambda}, \\ & \quad \sum_{j \in J(i)} \mu_j N_j y_{i,j} \geq \lambda_i, \quad i \in \mathcal{F}, \\ & \quad \sum_{i \in I(j)} y_{i,j} \leq 1, \quad j \in \mathcal{J}, \\ & \quad N \in \tilde{\mathcal{X}}^\lambda, \quad y_{i,j} \geq 0, \quad (i, j) \in E. \end{aligned} \quad (10)$$

We call a staffing vector determined through the solution of (10) a \wedge -based staffing. A standard argument shows that the optimization problem (10) can be solved by solving an associated LP that yields the same set of optimal solutions. Because we are not concerned with the way in which (10) is solved, we will use (10) directly. The next theorem provides the formal lower-bound result. Its proof is given in the online appendix.

THEOREM 2.1 (LOWER-BOUND CAPACITY). Consider the sequence of SBR systems and let $\{(N^\lambda, \pi^\lambda)\}$ be a sequence of asymptotically feasible staffing and routing rules such that

$$\liminf_{\lambda \rightarrow \infty} \frac{N_j^\lambda}{\lambda} > 0, \quad j \in \mathcal{J}.$$

Then, $\{N^\lambda, \lambda \geq 0\}$ satisfies that

$$\sum_{j \in \mathcal{J}} \mu_j N_j^\lambda \geq \lambda + \beta^* \sqrt{\lambda} + o(\sqrt{\lambda}), \quad (11)$$

where β^* is the \wedge -model parameter in (9).

Theorem 2.1 only provides a lower bound. We will next propose a solution that we will prove achieves this lower bound. The solution will be based on the \wedge -based staffing and the FQR routing rule.

3. The Proposed Solution

Our solution consists of a staffing component and a routing component. The staffing that we use is the \wedge -based staffing determined by an optimal solution to (10). For the routing component, we use FQR with ratios that will be explicitly determined as functions of the service-level targets T_i^λ .

Let $Q_i^\lambda(t)$ be the number of class- i customers in queue. Let $Z_{i,j}^\lambda(t)$ be the number of type- j servers busy giving service to class- i customers, so that $X_i^\lambda(t) := Q_i^\lambda(t) + \sum_{j \in \mathcal{J}} Z_{i,j}^\lambda(t)$ is the overall number of class- i customers present in the system at time t , and $I_j^\lambda(t) := N_j^\lambda - \sum_{i=1}^I Z_{i,j}^\lambda(t)$ be the number of idle agents in pool j at time t in the λ th system. Accordingly, $I_\Sigma^\lambda(t) := \sum_{j=1}^J I_j^\lambda(t)$ is the total number of idle agents in the system. Let $X_\Sigma^\lambda(t)$ be the overall number of customers in the system (in service and in queue), i.e.,

$$X_\Sigma^\lambda(t) := \sum_{i=1}^I X_i^\lambda(t) = \sum_{i=1}^I \left(Q_i^\lambda(t) + \sum_{j=1}^J Z_{i,j}^\lambda(t) \right),$$

and let $N_\Sigma^\lambda(t) := \sum_{j \in \mathcal{J}} N_j^\lambda$ be the aggregate number of agents. Below we use $\arg \max$ and let it have the standard definition; i.e., given a function $f: A \mapsto \mathbb{R}$, with A a finite set, let $\arg \max f := \{y \in A: f(y) = \max_{x \in A} f(x)\}$.

DEFINITION 3.1 (FQR FOR THE SBR MODEL). Given two probability vectors $v := \{v_j: j \in \mathcal{J}\}$ and $p := \{p_i: i \in \mathcal{F}\}$, FQR for the SBR model is defined as follows:

- Upon arrival of a class- i customer at time t , the customer will be routed to an available agent in pool j^* , where

$$j^* \equiv j^*(t) \in \arg \max_{j \in J(i), I_j^\lambda(t) > 0} \{I_j^\lambda(t) - v_j [X_\Sigma^\lambda(t) - N_\Sigma^\lambda]^-\};$$

i.e., the customer will be routed to an agent pool with the greatest idleness imbalance. If there are no such agents, the customer waits in queue i , to be served in order of arrival.

- Upon service completion by a type- j agent at time t , the agent will admit to service the customer from the head of queue i^* where

$$i^* \equiv i^*(t) \in \arg \max_{i \in I(j), Q_i^\lambda(t) > 0} \{Q_i^\lambda(t) - p_i [X_\Sigma^\lambda(t) - N_\Sigma^\lambda]^+\};$$

i.e., the agent will admit a customer from the queue with the greatest queue imbalance. If there are no such customers, the agent will remain idle.

Ties are broken in an arbitrary but consistent manner, so that the vector-valued stochastic process (Q^λ, Z^λ) is a continuous-time Markov chain (CTMC) with stationary transition probabilities.

To explicitly express the dependence on the vectors p and v , we will use the notation $FQR(p, v)$. We point out that if $p_i > 0$ for all $i \in \mathcal{J}$, then FQR is equivalently given by having each newly available agent choose the customer from the head of queue i^* , where

$$i^* \equiv i^*(t) \in \arg \max_{i \in I(j), Q_i^\lambda(t) > 0} \left\{ \frac{Q_i^\lambda(t)}{p_i} \right\}, \quad (12)$$

which makes the use of $[X_\Sigma^\lambda(t) - N_\Sigma^\lambda]^+$ unnecessary.

Choosing p and v . For the routing component of our solution, we will be using FQR with the ratio vectors (p^*, v^*) , where $v^* = (0, \dots, 0, 1)$ and p^* is the unique solution to

$$\mathbf{P}_{\mu_1}(\beta^*) e^{-(\lambda_{I-1}/\lambda p_{I-1})\beta^* \sqrt{\lambda} T_{I-1}^\lambda} = \alpha \quad \text{and} \quad (13)$$

$$\frac{p_i}{p_{I-1}} = \frac{\lambda_i T_i^\lambda}{\lambda_{I-1} T_{I-1}^\lambda},$$

for $\mathbf{P}_{\mu_1}(\beta)$ defined in (9). Because $\sqrt{\lambda} T_i^\lambda = \bar{T}_i$ (see Assumption 2.1) and $\lambda_i/\lambda_{I-1} = a_i/a_{I-1}$, the value of p_{I-1}^* is independent of λ . Consequently, so are the values p_i^* for $i \neq I-1$. The choice of the ratio vector p^* will be justified by (i) a sample-path version of Little's law that holds for the many-server service system, (ii) the SSC that is induced by FQR, and (iii) the fact that it performs asymptotically as efficiently as the \wedge model. Informally, SSC justifies the following sequence of approximations

$$\begin{aligned} P\{W_{I-1} > T_{I-1}\} &\approx P\{Q_{I-1} > \lambda_{I-1} T_{I-1}^\lambda\} \\ &\approx P\{p_{I-1} Q_\Sigma > \lambda_{I-1} T_{I-1}^\lambda\} \\ &\approx P\left\{Q_\Sigma^\lambda > \frac{\lambda_{I-1} T_{I-1}^\lambda}{p_{I-1}}\right\} \\ &\approx P\{Q_\Sigma^\lambda > 0\} e^{((\sum_{j \in \mathcal{J}} \mu_j N_j - \lambda)/\lambda)(\lambda_{I-1} T_{I-1}^\lambda/p_{I-1})}, \end{aligned}$$

where Q_Σ^λ is the steady-state queue length in the \wedge model that is constructed from the SBR system, as in the beginning of §2.2. The last step follows from simple expressions for the distribution of the queue length for the \wedge model. By the analysis of \wedge model in Armony (2005), the probability $P\{Q_\Sigma^\lambda > 0\}$ converges to $\mathbf{P}_{\mu_1}(\beta^*)$ if we use the \wedge -based staffing. Also, $\sum_{j \in \mathcal{J}} \mu_j N_j = \lambda + \beta^* \sqrt{\lambda} + o(\sqrt{\lambda})$. Hence,

$$\begin{aligned} P\{W_{I-1} > T_{I-1}\} &\approx P\left\{Q_\Sigma^\lambda > \frac{\lambda_{I-1} T_{I-1}^\lambda}{p_{I-1}}\right\} \\ &\approx \mathbf{P}_{\mu_1}(\beta^*) e^{-(\lambda_{I-1}/\lambda p_{I-1})\beta^* \sqrt{\lambda} T_{I-1}^\lambda}. \end{aligned} \quad (14)$$

Similar informal arguments can be repeated for each of the customer classes. Finally, because β^* was chosen so that the right-hand side in (14) equals α , the \wedge -based staffing and FQR should provide an asymptotically feasible solution to (3). Because the \wedge -based staffing is a lower bound, this solution is also asymptotically optimal. This informal argument is formalized in the next section.

4. Asymptotic Feasibility and Optimality

We begin to consider the limiting behavior as $\lambda \rightarrow \infty$. We will show that the \wedge -based staffing and FQR, with appropriately chosen ratios, yield an asymptotically feasible and optimal solution for (3). First, however, we consider the design of the system. In order to identify the design, it suffices to look at (10) with one constraint removed, i.e., consider the following nonlinear optimization problem:

$$\begin{aligned} &\text{Minimize } \sum_{j \in \mathcal{J}} c_j N_j \\ &\text{Subject to: } \sum_{j \in J(i)} \mu_j N_j y_{i,j} \geq \lambda_i, \quad i \in \mathcal{J}, \\ &\quad \sum_{i \in I(j)} y_{i,j} \leq 1, \quad j \in \mathcal{J}, \\ &\quad N \in \tilde{\mathcal{S}}^\lambda, \quad y_{i,j} \geq 0, (i, j) \in E. \end{aligned} \quad (15)$$

The solution to the mathematical program (15) can be regarded as a first-order deterministic fluid approximation for the SBR system, as in Whitt (2006). From that point of view, given a selected solution (\bar{N}, \bar{y}) , we would then use \bar{N} to provide an initial estimate of the staffing and \bar{y} to provide an initial estimate of the appropriate routing. We point out that the solution to (15) is independent of λ . Indeed, by the definition of $\tilde{\mathcal{S}}^\lambda$ and the assumption that $\lambda_i = a_i \lambda$, (15) is equivalent to the mathematical program:

$$\begin{aligned} &\text{Minimize } \sum_{j \in \mathcal{J}} c_j \nu_j \\ &\text{Subject to: } \sum_{j \in J(i)} \mu_j \nu_j x_{i,j} \geq a_i, \quad i \in \mathcal{J}, \\ &\quad \sum_{i \in I(j)} x_{i,j} \leq 1, \quad j \in \mathcal{J}, \\ &\quad A \nu \leq b, \\ &\quad \nu_j \geq 0, \quad x_{i,j} \geq 0, j \in \mathcal{J}, (i, j) \in E. \end{aligned} \quad (16)$$

Both mathematical programs (15) and (16) can be replaced with linear programs (LP) that yield the same optimal solution. Henceforth, we only refer to optimal solutions of (15), without considering how they are obtained. We denote an optimal solution to (15) by $(\bar{N}^\lambda, \bar{y}^\lambda)$. Our first assumption requires a weak form of uniqueness of optimal solutions to (15).

ASSUMPTION 4.1 (UNIQUENESS OF STAFFING). Fix λ and let $(\bar{N}^\lambda, \bar{y}^\lambda)$ and $(\tilde{N}^\lambda, \tilde{y}^\lambda)$ be two optimal solutions to (15). Then $\bar{N}^\lambda = \tilde{N}^\lambda$.

We note that due to the equivalence between (15) and (16), if Assumption 4.1 holds for a given λ , then it holds for all λ . Similarly, the equivalence between (15) and (16) implies that if $\bar{N}_j^\lambda > 0$ for one λ , then the same holds for all values of λ . Informally, Assumption 4.1 is required because we will want to use the optimal solutions of (15) as a first-order (fluid) approximation for the staffing (and not only

the staffing cost) that we get from the \wedge -based staffing as defined by the solution to (10). In some simplified settings such as the one considered in §5.4, this assumption can be removed.

The second assumption is a critical loading assumption, needed to put the system in heavy traffic.

ASSUMPTION 4.2 (CRITICAL LOADING). For any $\lambda \geq 0$, $\sum_{j \in \mathcal{J}} \mu_j \bar{N}_j^\lambda = \lambda$ for any optimal solution $(\bar{N}^\lambda, \bar{y}^\lambda)$ to (15).

Finally, we make the following structural assumption. Below, E and V are as defined in the beginning of §2. Also, we say that a graph is connected if there exists a path between every two nodes in the graph.

ASSUMPTION 4.3 (CONNECTED ROUTING GRAPH). For any $\lambda \geq 0$, there exists an optimal solution $(\bar{N}^\lambda, \bar{y}^\lambda)$ for (15) such that the graph $\mathcal{G}(\bar{N}^\lambda, \bar{y}^\lambda) := \{(i, j) \in \mathcal{J} \times \mathcal{J} : \bar{y}_{i,j}^\lambda > 0\}$ is a connected subgraph of $G(V, E)$.

As before, the equivalence between (15) and (16) guarantees that if Assumption 4.3 holds for a given λ , then it holds for all λ . This connected-graph assumption is crucial for the ability to instantaneously balance the system by re-directing capacity from one customer class to the other; see §2.7 of Atar (2005) for elaboration. Assumptions 4.1–4.3 are assumed to hold throughout the rest of the paper. With the above definitions, we can state our asymptotic optimality result.

THEOREM 4.1 (ASYMPTOTIC OPTIMALITY FOR THE SBR MODEL WITH POOL-DEPENDENT RATES). Suppose that any optimal solution for (15) has $\bar{N}_j^\lambda > 0$ for all $j \in \mathcal{J}$. Let N^λ be determined through the \wedge -based staffing in (10) with β^* as in (9). Set π^λ to FQR(p^*, v^*) with p^* as in (13) and $v^* = (0, \dots, 0, 1)$. Then, the sequence $\{(N^\lambda, \pi^\lambda)\}$ is asymptotically optimal for (3).

REMARK 4.1 (CHOOSING THE RATIO VECTOR v^*). In light of our SSC result in Theorem 3.1 of Gurvich and Whitt (2009a), the choice $v^* = (0, 0, \dots, 1)$ will cause all the idleness to be concentrated in pool J , which is the slowest agent-pool. This choice guarantees that all the faster servers will be constantly busy, thus maximizing the depletion rate of customers from the system. Informally, then, this choice of v^* minimizes the aggregate queue length in the system by maximizing the depletion rate. Because this observation holds for any staffing level, this choice of v^* is essential for the minimization of the number of agents required to achieve the aggregate waiting-time constraints. Once the aggregate queue length is minimized, it only remains to distribute it in a proper way to ensure that the SL constraints are met. The queue-ratio vector, p^* , takes care of this task. \square

Theorem 4.1 illustrates one of the key benefits of FQR. Although the \wedge model is a more efficient system, FQR allows the SBR system to work as efficiently, asymptotically, making the staffing of the \wedge model sufficient also for the SBR system.

Asymptotic Feasibility for (2). We now discuss an asymptotically feasible solution for the SBR problem (2). Although this formulation is somewhat problematic, as discussed in §2, it is very common in industry. Hence, it is of interest to discuss the construction of feasible solutions for this problem. We now define p^* to be

$$p_i^* := \frac{\lambda_i T_i^\lambda}{\sum_{k \in \mathcal{J}} \lambda_k T_k^\lambda} = \frac{a_i \bar{T}_i}{\sum_{k \in \mathcal{J}} a_k \bar{T}_k}, \quad (17)$$

and redefine β^* to be the unique solution of

$$\mathbf{P}_{\mu_1}(\beta) e^{-\beta \sqrt{\lambda} \sum_{i \in \mathcal{J}} (\lambda_i / \lambda) T_i^\lambda} = \alpha, \quad (18)$$

where $\mathbf{P}_{\mu_1}(\beta)$ defined in (9). As before, we observe that the vector p^* is independent of λ , because $\lambda_i / \lambda = a_i$ and Assumption 2.1 holds.

THEOREM 4.2 (ASYMPTOTIC FEASIBILITY FOR THE SBR MODEL WITH POOL-DEPENDENT RATES). Suppose that any optimal solution for (15) has $\bar{N}_j^\lambda > 0$ for all $j \in \mathcal{J}$. Let N^λ be determined through the \wedge -based staffing in (10) with β^* as in (18). Set π^λ to FQR(p^*, v^*) with p^* as in (17) and $v^* = (0, \dots, 0, 1)$. Then, the sequence $\{(N^\lambda, \pi^\lambda)\}$ is asymptotically feasible for (2).

Both Theorems 4.1 and 4.2 rely on the fact that FQR admits an important SSC result by which, asymptotically, the queues of class i are equal to the proportion p_i^* of the aggregate queue length, and the number of idle servers in pool j is equal to a proportion v_j^* of the aggregate number of idle servers. Somewhat informally, the SSC results guarantee that with the \wedge staffing and FQR we will have that, for all $i \in \mathcal{J}$ and $j \in \mathcal{J}$

$$\frac{Q_i^\lambda(t)}{\sqrt{\lambda}} - v_j \frac{Q_\Sigma^\lambda(t)}{\sqrt{\lambda}} \Rightarrow 0, \quad \text{as } \lambda \rightarrow \infty, \quad \text{and}$$

$$\frac{I_j^\lambda(t)}{\sqrt{\lambda}} - p_i \frac{I_\Sigma^\lambda(t)}{\sqrt{\lambda}} \Rightarrow 0, \quad \text{as } \lambda \rightarrow \infty,$$

where $Q_\Sigma^\lambda(t)$ and $I_\Sigma^\lambda(t)$ are, respectively, the aggregate queue length and aggregate number of idle agents at time t . The SSC result for this setting is a corollary of our more general result in Theorem 3.1 of Gurvich and Whitt (2009a).

5. Other Formulations

This section is dedicated to alternative formulation of the call-center optimization problem.

5.1. Constraints on Average Delay

Here, we consider the formulation.

$$\begin{aligned} & \text{minimize} && \sum_{j \in \mathcal{J}} c_j N_j \\ & \text{subject to} && \limsup_{T \rightarrow \infty} \bar{W}_i^{\lambda, T} \leq T_i^\lambda, \quad i \in \mathcal{J}, \\ & && N \in \mathcal{N}^\lambda, \quad \pi \in \Pi, \end{aligned} \quad (19)$$

where

$$\bar{W}_i^{\lambda, T} := \frac{\sum_{k=1}^{A_i^\lambda(T)} w_{i,k}^\lambda}{A_i^\lambda(T)}.$$

DEFINITION 5.1 (ASYMPTOTIC FEASIBILITY FOR (19)). A sequence of staffing vectors and routing policies $\{(N^\lambda, \pi^\lambda)\}$ is asymptotically feasible for (19) if: (a) $\pi^\lambda \in \Pi$, for all λ , and (b) for every $\epsilon > 0$, there exists $T^*(\epsilon)$ such that, for all $T \geq T^*(\epsilon)$,

$$\limsup_{\lambda \rightarrow \infty} P \left\{ \frac{W_i^{\lambda, T}}{T_i^\lambda} \geq 1 + \epsilon \right\} \leq \epsilon, \quad i \in \mathcal{F}. \quad (20)$$

The proposed solution in this case is as follows:

• **Staffing:** We use the optimal solution to (10) with β^* now given by the unique solution to

$$\frac{\mathbf{P}_{\mu_1}(\beta)}{\beta} = \sum_{i \in \mathcal{F}} a_i \sqrt{\lambda} T_i^\lambda = \sum_{i \in \mathcal{F}} a_i \bar{T}_i, \quad (21)$$

where again $\mathbf{P}_{\mu_1}(\beta)$ is defined in (9).

• **Routing:** FQR with ratio vectors $v^* = (0, 0, \dots, 1)$ and p^* that is given by

$$p_i^* := \frac{\lambda_i T_i^\lambda}{\sum_{k \in \mathcal{F}} \lambda_k T_k^\lambda} = \frac{a_i \bar{T}_i}{\sum_{k \in \mathcal{F}} a_k \bar{T}_k}. \quad (22)$$

THEOREM 5.1 (ASYMPTOTIC OPTIMALITY FOR THE AVERAGE-WAITING-TIME FORMULATION). Suppose that any optimal solution for (15) has $\bar{N}_j^\lambda > 0$ for all $j \in \mathcal{F}$. Let N^λ be determined through the \wedge -based staffing in (10) with β^* as in (21). Set π^λ to FQR(p^*, v^*), where p^* is as in (22) and $v^* = (0, \dots, 0, 1)$. Then, the sequence $\{(N^\lambda, \pi^\lambda)\}$ is asymptotically optimal for (19).

5.2. Adding Customer Abandonment

In this section we augment our model by introducing customer abandonment. Specifically, we assume that a class- i customer has an exponential patience with rate θ_i . If his patience expires before he is admitted to service, the customer will abandon. Patience times of different customers are mutually independent. We will formulate an optimization problem for the abandonment model and provide the corresponding asymptotic feasibility and optimality results. Our asymptotic optimality result for this augmented model is limited to the case in which all customer classes have the same abandonment rate, i.e., when $\theta_i \equiv \theta$. It is of practical interest, however, that we can provide asymptotically feasible solutions for the case in which these rates are different.

To prove the asymptotic results for the abandonment model, we exploit Armony and Mandelbaum (2008), which extends the results of Armony (2005) to the \wedge model with customer abandonments. Using Armony and Mandelbaum (2008), we are able to provide proofs that parallel those for the nonabandonment case, repeating the analyses of (2), (3),

and (19). Asymptotic feasibility results can be obtained for general patience rates θ_i , whereas asymptotic optimality can be obtained only for the homogeneous case.

Rather than repeating the analysis for formulations (2) and (3), we introduce and solve a different formulation with constraint on the fraction of abandoning customers. To formulate this problem, let $L_i^\lambda(t)$ be the number of class- i customers that abandoned before being served, up to time t . The fraction of customers that abandoned among those that were initially in the queue (at time $t = 0$) and those that arrived after time 0, is then given by

$$\text{Ab}_i^{\lambda, T} := \frac{L_i^\lambda(T)}{Q_i^\lambda(0) + A_i^\lambda(T)}. \quad (23)$$

We then consider the formulation

$$\begin{aligned} & \text{minimize} \quad \sum_{j \in \mathcal{F}} c_j N_j \\ & \text{subject to} \quad \limsup_{T \rightarrow \infty} \text{Ab}_i^{\lambda, T} \leq \alpha_i^\lambda, \quad i \in \mathcal{F}, \\ & \quad \quad \quad N \in \mathcal{A}^\lambda, \quad \pi \in \Pi, \end{aligned} \quad (24)$$

where we assume that $\alpha_i^\lambda = \bar{\alpha}_i / \sqrt{\lambda}$ for some strictly positive constants $\{\bar{\alpha}_i, i \in \mathcal{F}\}$, in order to place the system in the QED regime.

DEFINITION 5.2 (ASYMPTOTIC FEASIBILITY FOR (24)). A sequence of staffing vectors and routing policies $\{(N^\lambda, \pi^\lambda)\}$ is asymptotically feasible for (24) if: (a) $\pi^\lambda \in \Pi$ for all λ ; and (b) for every $\epsilon > 0$, there exists $T^*(\epsilon)$ such that for all $T \geq T^*(\epsilon)$,

$$\limsup_{\lambda \rightarrow \infty} P \left\{ \frac{\text{Ab}_i^{\lambda, T}}{\alpha_i^\lambda} \geq 1 + \epsilon \right\} \leq \epsilon, \quad i \in \mathcal{F}. \quad (25)$$

We propose the following staffing and routing solution:

• **Staffing:** Use the optimal solution to (10) with β^* now given by the unique solution to

$$\begin{aligned} \bar{\alpha} &= \sqrt{\bar{\theta}} \mathbf{P}_{\mu_1, \bar{\theta}}(\beta) \left[h\left(\frac{\beta}{\sqrt{\bar{\theta}}}\right) - \frac{\beta}{\sqrt{\bar{\theta}}} \right] \quad \text{and} \\ \mathbf{P}_{\mu_1, \bar{\theta}}(\beta) &:= \left[1 + \frac{\sqrt{\bar{\theta}} h(\beta/\sqrt{\bar{\theta}})}{\sqrt{\mu_1} h(-\beta/\sqrt{\mu_1})} \right]^{-1}, \end{aligned} \quad (26)$$

where $h(\cdot) := \phi(\cdot)/(1 - \Phi(\cdot))$.

$$\bar{\theta} := \sum_{i \in \mathcal{F}} p_i \theta_i \quad \text{and} \quad \bar{\alpha} = \sum_{i \in \mathcal{F}} a_i \alpha_i. \quad (27)$$

• **Routing:** FQR with ratio vectors $v^* = (0, 0, \dots, 1)$ and p^* given by

$$p_i^* := \frac{\lambda_i \alpha_i^\lambda / \theta_i}{\sum_{k \in \mathcal{F}} \lambda_k \alpha_k^\lambda / \theta_k} = \frac{a_i \bar{\alpha}_i / \theta_i}{\sum_{k \in \mathcal{F}} a_k \bar{\alpha}_k / \theta_k}, \quad (28)$$

In (26), $\mathbf{P}_{\mu_1, \bar{\theta}}(\beta)$ is the asymptotic delay probability in the corresponding \wedge model under the FSF policy, with the patience rate θ and having $\sum_{j \in \mathcal{F}} \mu_j N_j = \lambda + \beta \sqrt{\lambda} + o(\sqrt{\lambda})$; see Propositions 4.5 and 4.6 in Armony and Mandelbaum (2008).

THEOREM 5.2 (ASYMPTOTIC FEASIBILITY AND OPTIMALITY FOR THE ABANDONMENT FORMULATION). *Suppose that any optimal solution for (15) has $\bar{N}_j^\lambda > 0$ for all $j \in \mathcal{J}$. Let N^λ be determined through the \wedge -based staffing in (10) with β^* in (26). Set π^λ to FQR(p^*, v^*), where p^* is as in (28) and $v^* = (0, \dots, 0, 1)$. Then, the sequence $\{(N^\lambda, \pi^\lambda)\}$ is asymptotically feasible for (24). If, in addition, $\theta_i \equiv \theta$ for all i , then the sequence $\{(N^\lambda, \pi^\lambda)\}$ is also asymptotically optimal.*

We prove Theorem 5.2 in the e-companion. The required argument is similar to the previous case without abandonments. The key step is to show that with FQR, the SBR model with abandonment is asymptotically equivalent (in terms of the aggregate number of customers in system) to a \wedge model in which the customers' patience is exponential with a rate that is averaged using the ratio vector p of FQR in Equation (27). For homogeneous patience rates, namely when $\theta_i \equiv \theta$, we show that a \wedge -based staffing (modified for the abandonment case) provides a lower bound on the staffing costs. Asymptotic optimality then follows from the asymptotic feasibility.

EXAMPLE 1 (A TWO-CLASS TWO-POOL SYSTEM). We apply the proposed staffing-and-routing solution to a two-class two-pool system. We assume that pool-1 servers can serve only class-1 customers, whereas pool-2 servers are cross-trained. The resulting N-model is depicted in Figure 2. The customer arrival and patience parameters are, respectively, $(\lambda_1, \lambda_2) = (100, 50)$ and $(\theta_1, \theta_2) = (2, 1)$. Because we are considering a fixed system, we omit the superscript λ from all notation. Because we are considering a setting with non-homogeneous patience rates, we aim only to show the feasibility of our solution.

To complete the model description, let the (pool-dependent) service rates be $(\mu_1, \mu_2) = (1.5, 1)$; let $c_1 = c_2 = c$; and let $\mathcal{A} = \{N \in \mathbb{Z}_+^2 : N_1 \leq 50\}$. In particular, the

number of agents in the first pool can be at most 50. Finally, we assume that the abandonment constraints are 3% for class 1 and 5% for class 2, i.e., $(\alpha_1, \alpha_2) = (0.03, 0.05)$.

With these parameters, we construct our (asymptotically feasible) \wedge -based staffing and FQR routing solution. The parameters for FQR are $p_1^* = 0.375$, $p_2^* = 0.625$, $\bar{\theta} = 1.375$, and $\bar{\alpha} = 0.3266$. From (26), for the \wedge -based staffing we have $\beta^* = 0.03926$, so that we need $\mu_1 N_1 + \mu_2 N_2 = 150 + \beta^* \sqrt{150} \geq 125.48$. Accordingly, we set $N_1 = 50$ and $N_2 = \lceil (125.4808 - 50)/\mu_2 \rceil = 76$.

We simulate this N model with the specified staffing and FQR(p^*, v^*). We run 3,000 replications of the system, each up to $T = 500$. The graph in Figure 2 displays the average proportion of abandoning customers for each customer class and for each time unit (as a function of time). Evidently, the proportion of abandonments is below the target (for class 1) or only slightly above (for class 2). Also, we find that

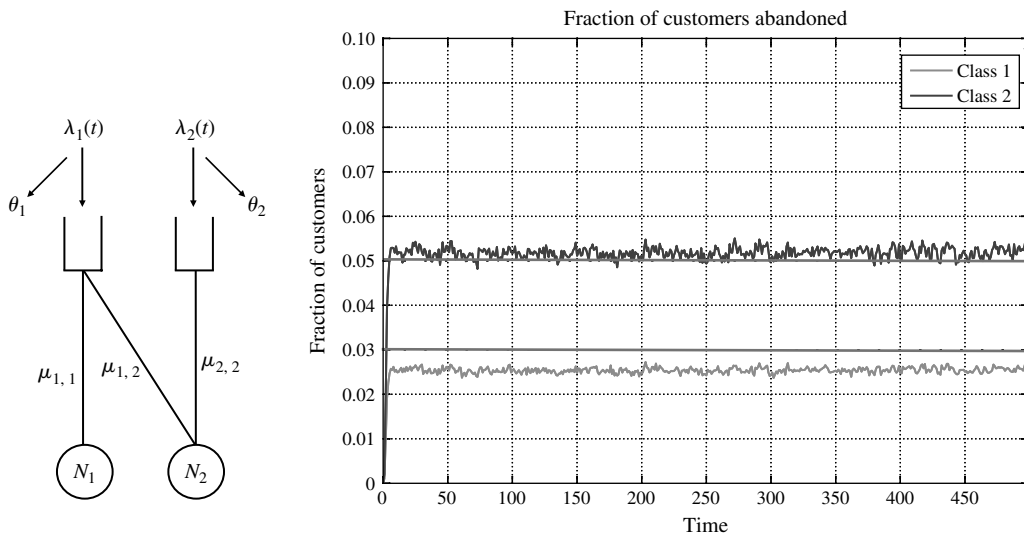
$$P \left\{ \frac{\text{Ab}_i^T}{\alpha_i} \geq 1.02 \right\} \leq 0.02, \quad i = 1, 2, \tag{29}$$

where $P\{\cdot\}$ should be interpreted here as the empirical probability distribution over the 3,000 replications, whereas (29) corresponds to the asymptotic feasibility condition (25) in Definition 5.2 with ϵ there taken to be 0.02. \square

5.3. Designated-Service Constraints

In practice, it is natural to require that most customers receive their *designated service*, i.e., service by the type of agent that the system designates for them. A good example is a multilingual call center, where one would prefer that Spanish-speaking customers be served by agents whose dominant language is Spanish, French-speaking customers be served by agents whose dominant language is French, and so forth. In other settings, the interpretation of

Figure 2. A two-class two-pool N model and the simulation results.



Copyright: INFORMS holds copyright to this *Articles in Advance* version, which is made available to institutional subscribers. The file may not be posted on any other website, including the author's site. Please send any questions regarding this policy to permissions@informs.org.

designated service might be different. To treat designated-service constraints, we now assume that for each customer class i there is one agent type designated to provide service to that class; the agent type designated for class i is $j(\{i\})$; that pool of agents can only serve class i .

To impose a lower-bound constraint on the proportion of customers that receive their designated service, we let $D_i^{\lambda, T}$ be the proportion of the arriving class- i customers that are routed to their designated agents (the agents in pool $j(\{i\})$) by time T . Letting $\delta_i > 0$ be the nondesignated-service proportion upper bound, the constraints are formulated as

$$\liminf_{T \rightarrow \infty} D_i^{\lambda, T} \geq 1 - \delta_i, \quad i \in \mathcal{J}.$$

The optimization problem that we consider is then given by

$$\begin{aligned} & \text{minimize} && \sum_{j \in \mathcal{J}} c_j N_j \\ & \text{subject to} && \limsup_{T \rightarrow \infty} \bar{W}^{\lambda, T} \leq T_I^\lambda, \\ & && \limsup_{T \rightarrow \infty} \bar{F}_i^{\lambda, T}(T_i^\lambda) \leq \alpha, \quad 1 \leq i \leq I - 1, \quad (30) \\ & && \liminf_{T \rightarrow \infty} D_i^{\lambda, T} \geq 1 - \delta_i, \quad i \in \mathcal{J}, \\ & && N \in \mathcal{A}^\lambda, \quad \pi \in \Pi. \end{aligned}$$

We now show that the designated service constraints can be incorporated into the framework of §4 by appropriately redefining the set \mathcal{A}^λ . To this end, we extend the definition of asymptotic feasibility as follows:

DEFINITION 5.3 (ASYMPTOTIC FEASIBILITY WITH DESIGNATED-SERVICE CONSTRAINTS). A sequence of $\{(N^\lambda, \pi^\lambda)\}$ is asymptotically feasible for (30) if (a) $\pi^\lambda \in \Pi$, for all λ , and (b) for every $\epsilon > 0$, there exists $T^*(\epsilon)$ such that, for all $T \geq T^*(\epsilon)$, Equations (4) and (5) hold, as well as

$$\limsup_{\lambda \rightarrow \infty} P\{D_i^{\lambda, T} \leq 1 - \delta_i - \epsilon\} \leq \epsilon, \quad i \in \mathcal{J}. \quad (31)$$

The follow lemma provides a property that all asymptotically feasible solutions must satisfy.

LEMMA 5.1. *If $\{(N^\lambda, \pi^\lambda)\}$ is a sequence of asymptotically feasible staffing and routing rules in the sense of Definition 5.3, then*

$$\liminf_{\lambda \rightarrow \infty} \frac{\mu_{j(i)} N_{j(i)}^\lambda}{\lambda_i} \geq (1 - \delta_i), \quad i \in \mathcal{J}. \quad (32)$$

Lemma 5.1 suggests that one may incorporate the setting with designated-service constraints within the framework of §2 by replacing the optimization problem (30) with

$$\begin{aligned} & \text{minimize} && \sum_{j \in \mathcal{J}} c_j N_j \\ & \text{subject to} && \limsup_{T \rightarrow \infty} \bar{W}^{\lambda, T} \leq T_I^\lambda \\ & && \limsup_{T \rightarrow \infty} \bar{F}_i^{\lambda, T}(T_i^\lambda) \leq \alpha, \quad 1 \leq i \leq I - 1, \quad (33) \\ & && N \in \mathcal{C}^\lambda, \quad \pi \in \Pi, \end{aligned}$$

where

$$\begin{aligned} \mathcal{C}^\lambda &:= \mathcal{A}^\lambda \cap \mathcal{B}^\lambda \quad \text{and} \\ \mathcal{B}^\lambda &:= \left\{ N \in \mathbb{Z}_+^J : N_{j(i)} \geq (1 - \delta_i) \frac{\lambda_i}{\mu_{j(i)}}, i \in \mathcal{J} \right\}. \quad (34) \end{aligned}$$

Note that the set \mathcal{C}^λ fits in the framework of §2.1 as it is defined by linear constraints. In particular, (33) is a special case of (3) obtained by letting the set A^λ there be equal to \mathcal{C}^λ . The formal asymptotic feasibility result for this section appears in the next proposition.

THEOREM 5.3. *Suppose that the conditions of Theorem 4.1 hold with respect to formulation (33). Let N^λ be the asymptotically optimal \wedge -based staffing based on (10), with \mathcal{A}^λ replaced by \mathcal{C}^λ in (34) and set π^λ to FQR(p^*, v^*) with p^* as in (13) and $v^* = (0, \dots, 0, 1)$. Then, $\{(N^\lambda, \pi^\lambda)\}$ is asymptotically feasible for (30) in the sense of Definition 5.3. It is asymptotically optimal for (30) if one of the following holds: (i) $A^\lambda \supseteq B^\lambda$, or (ii) $c_j \equiv c$ and $\mu_j \equiv \mu$.*

We emphasize that the asymptotically optimal solution to (33), although asymptotically feasible for (30), need not be, in general, asymptotically optimal. The inequalities imposed by the set \mathcal{B}^λ might be too restrictive. Actually, as the proof of Proposition 5.3 reveals, we could have \mathcal{B}^λ defined through

$$\mathcal{B}^\lambda = \left\{ N \in \mathbb{Z}_+^J : N_{j(i)} \geq (1 - \delta_i) \frac{\lambda_i}{\mu_{j(i)}} - K\sqrt{\lambda}, i \in \mathcal{J} \right\}$$

for some $K > 0$. This would be less restrictive, but would suffice to generate an asymptotically feasible solution for (30).

In the next section we consider a setting in which $c_j \equiv c$ and $\mu_j \equiv \mu$. There, as in the second part of Proposition 5.3, we will be able to replace (30) with (33) without compromising asymptotic optimality.

5.4. Common Service Rates

This section is devoted to a simple setting: a common service rate μ , no abandonments, a common cost c for all agents, and $\mathcal{A}^\lambda = \mathbb{Z}_+$, as considered in Wallace and Whitt (2005). We have three purposes: first, to contrast the FQR-based solution with the simulation-based approach of Wallace and Whitt (2005), second, to illustrate an explicit construction of a system design when costs and system constraints do not pose significant restrictions; and third, to illustrate the diminishing-return property of flexibility. The optimization problem that we consider in this section is (30), but with $c_j \equiv c$, $j \in \mathcal{J}$, and $\mathcal{A}^\lambda \equiv \mathbb{Z}_+^J$.

Because we have a common service rate, we can staff using (nonasymptotic) formulas for the $M/M/N$ queue instead of the asymptotic expressions associated with the \wedge -based staffing in §4. (In this setting, these staffing methods are asymptotically equivalent.) To specify the staffing method here, let

$$N_\Sigma^\lambda := \min\{N \in \mathbb{Z}_+ : E[W_{\lambda, \mu}^{\text{FCFS}}] \leq \lambda T_I^\lambda\}, \quad (35)$$

where $W_{\lambda, \mu}^{\text{FCFS}}$ is the steady-state waiting time in an $M/M/N$ queue with arrival rate λ and service rate μ . For routing, we also can use nonasymptotic expressions. Specifically, we can use FQR with ratio vector p given by

$$P\{W_{\lambda, \mu}^{\text{FCFS}} > 0\} e^{-\mu(\bar{N}_{\Sigma} - \lambda/\mu)(\lambda_{I-1} T_{I-1}^{\lambda} / \lambda p_{I-1})} = \alpha \quad \text{and} \quad (36)$$

$$\frac{p_i}{p_{I-1}} = \frac{\lambda_i T_i^{\lambda}}{\lambda_{I-1} T_{I-1}^{\lambda}},$$

Note that the ratio vector p here does depend on λ .

Rather than assuming that the design is given, we allow ourselves in this section to choose the design. In doing this, we will take into account the designated-service constraint. We will incorporate this constraint from §5.3. The specific design we suggest here is based on a concatenation of M systems, which we call the *generalized M (GM) model*. An example of a GM model with three customer classes is depicted in Figure 3. The GM model has a routing graph constructed by allowing only edges of form $(i, j(\{i\}))$, $i \in \mathcal{I}$, and $(i, j(\{i, i+1\}))$, $i \in \mathcal{I} \setminus I$ (\mathcal{I} excluding the element I). The GM model is relatively inexpensive in terms of cross-training, because it uses agents with at most two skills, and only a limited number with two skills.

The solution we propose is as follows:

• **Design: generalized M model (GM).** Use a GM model.

• **Staffing: Single-Class Staffing (SCS).** Determine the overall number of agents, \bar{N}_{Σ} using (35). Then allocate agents to the pools by

$$- N_{j(\{i\})} = (1 - \delta/2)(\lambda_i/\lambda)\bar{N}_{\Sigma}, \quad i \in \{1, I\}, \quad N_{j(\{i\})} = (1 - \delta)(\lambda_i/\lambda)\bar{N}_{\Sigma} \quad \text{for all } i = 2, \dots, I-1, \quad \text{and}$$

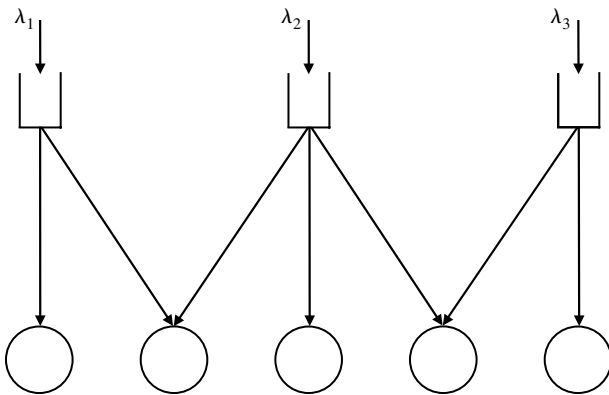
$$- N_{j(\{i, i+1\})} = (\delta(\lambda_i + \lambda_{i+1})/(2\lambda))\bar{N}_{\Sigma} \quad \text{for all } i = 1, \dots, I-1,$$

where $\delta := \min\{\delta_i: 1 \leq i \leq I\} > 0$.

• **Control: Fixed-Queue-Ratio (FQR).** Use FQR with p as defined in (36) and v defined by

$$v_{j(\{i, i+1\})} := \frac{1}{I-1} \quad \text{for all } i \in \mathcal{I} \setminus I. \quad (37)$$

Figure 3. The generalized M model for three classes.



The common service rate allows us to use any vector v in the control step above. This stands in contrast to the case of different service rates, where we needed to use $v^* = (0, \dots, 0, 1)$. The specific vector in (37) is designed to increase the amount of designated service by forcing the system to route customers that find agents idle in both pools $j(\{i\})$ and $j(\{i, i+1\})$ to the designated agents in pool $j(\{i\})$. One could also modify FQR so that all customers that find agents idle in more than one agent pool that can serve them will go to the designated agent pool $j(\{i\})$. This modification is guaranteed to achieve, asymptotically, the same performance as FQR. Using the results in §4, the above combined design-staffing-and-control solution can be shown to be asymptotically optimal as the arrival rate grows:

THEOREM 5.4 (ASYMPTOTIC OPTIMALITY FOR THE SBR MODEL WITH COMMON SERVICE RATES). *Consider the simple SBR model specified above and assume that the GM design is used. Let N^{λ} be determined by SCS staffing and set π^{λ} to FQR(p^* , v^*) with p^* as in (36) and v^* as in (37). Then, the sequence $\{(N^{\lambda}, \pi^{\lambda})\}$ is asymptotically optimal for (30) with $c_j \equiv c$ and $\mathcal{A}^{\lambda} \equiv \mathbb{Z}_+^J$.*

6. Conclusions

In this paper we have proposed the *fixed-queue-ratio* (FQR) routing scheme for the real-time routing of customers in call centers with multiple customer classes and multiple agent types operating under QoS constraints. FQR routing facilitates the construction of *combined staffing-design-and-routing solutions* for some settings of the complicated *skill-based-routing* (SBR) problem, with precisely specified goals. In this paper, we used FQR to construct an asymptotically optimal solution for the staffing-and-routing problem subject to QoS constraints. The key assumption that we made is that the service rates are pool dependent. However, as discussed in §2, this is not enough, and we need to be careful about the formulation; to get asymptotic optimality, we need to replace the initial formulation (2) with the best-effort formulation in (3); Theorem 4.1 shows that FQR is asymptotically optimal in this setting. FQR also produces asymptotic optimality for other important formulations, specified in §5. Some modification is needed for each new formulation, but a version of FQR applies in each case. A key component of the proof was showing that our SBR problem is asymptotically equivalent to the \wedge model previously analyzed by Armony (2005).

It is especially instructive to see what can be done in the special case of a common service rate μ and a common agent cost c considered in §5.4, which was previously considered by Wallace and Whitt (2005) using an iterative simulation-based staffing algorithm. In that case, we need neither the \wedge -based staffing in (10), nor the optimization problem (15). Consequently, for this special case we do not need Assumptions 4.1–4.3. This simple model illustrates how FQR simplifies tremendously the construction of the

joint design-staffing-and-control solution, by allowing one to ignore the SL constraints when making the design and staffing decisions. The FQR routing will take care of those through a simple choice of the ratio vector. This essential decoupling of the design, staffing and control decisions is beneficial for applications, because in practice the design and staffing decisions are indeed often made in advance, and cannot easily be adjusted in real time. This stands in contrast to Wallace and Whitt (2005), where the numbers of agents in the service pools need to be fine-tuned through simulation to meet the SL constraints.

Finally, it is significant that the GM design used in §5.4 uses only limited flexibility. In particular, it uses agents that have at most two skills, and then only a limited number with two skills. Still, with this limited amount of flexibility, the SBR system performs, asymptotically, as efficiently as the single-class single-pool $M/M/N$ queue. Although this idea was communicated in Wallace and Whitt (2005), no mathematical results were established there.

7. Electronic Companion

An electronic companion to this paper is available as part of the online version that can be found at <http://or.journal.informs.org/>.

Acknowledgments

This research is based on the first author's doctoral dissertation at Columbia University. The authors are grateful to Avi Mandelbaum and Mor Armony for fruitful discussions and to Zohar Feldman for contributions to the simulation. The second author was supported by NSF grant DMI-0457095.

References

Armony, M. 2005. Dynamic routing in large-scale service systems with heterogeneous servers. *Queueing Systems* **51**(3–4) 287–329.
Armony, M., A. Mandelbaum. 2008. Routing and staffing in large-scale service systems: The case of homogeneous impatient customers and heterogeneous servers. Working paper, New York University, New York,

and Technion—Israel Institute of Technology, Haifa, Israel.
Atar, R. 2005. Scheduling control for queueing systems with many servers: Asymptotic optimality in heavy traffic. *Ann. Appl. Probab.* **15**(4) 2606–2650.
Bassamboo, A., A. Zeevi. 2008. Staffing telephone call centers subject to service-level constraints: An approximate approach via constraint dualization. Working paper, Northwestern University, Evanston, IL, and Columbia University, New York.
Bassamboo, A., J. M. Harrison, A. Zeevi. 2006a. Dynamic routing and admission control in high volume service systems: Asymptotic analysis via multi-scale fluid limits. *Queueing Systems* **51**(3–4) 249–285.
Bassamboo, A., J. M. Harrison, A. Zeevi. 2006b. Design and control of a large call center: Asymptotic analysis of an LP-based method. *Oper. Res.* **54**(3) 419–435.
Borst, S., A. Mandelbaum, M. Reiman. 2004. Dimensioning large call centers. *Oper. Res.* **52**(1) 17–34.
Feldman, Z., I. Gurvich, W. Whitt. 2007. Managing quality of service in call centers via queue-ratio routing: Asymptotic analysis and simulation-based optimization. Working paper, Columbia University, New York.
Gans, N., G. Koole, A. Mandelbaum. 2003. Telephone call centers: Tutorial, review and research prospects. *Manufacturing Service Oper. Management* **5**(2) 79–141.
Gurvich, I., W. Whitt. 2009a. Queue-and-idleness-ratio controls in many-server service systems. *Math. Oper. Res.* **34**(2) 363–396.
Gurvich, I., W. Whitt. 2009b. Scheduling flexible servers with convex delay costs in many-server service systems. *Manufacturing Service Oper. Management* **11**(2) 237–253.
Gurvich, I., M. Armony, A. Mandelbaum. 2008. Service-level differentiation in call centers with fully flexible servers. *Management Sci.* **54**(2) 279–294.
Gurvich, I., J. Luedtke, T. Tezcan. 2008. Staffing call centers with uncertain demand forecasts: A chance constrained optimization approach. Working paper, Northwestern University, Evanston, IL.
Halfin, S., W. Whitt. 1981. Heavy-traffic limits for queues with many exponential servers. *Oper. Res.* **29**(3) 567–587.
Mandelbaum, A., A. Stolyar. 2004. Scheduling flexible servers with convex delay costs: Heavy-traffic optimality of the generalized $c\mu$ -rule. *Oper. Res.* **52**(6) 836–855.
Van Mieghem, J. A. 1995. Dynamic scheduling with convex delay costs: The generalized $c\mu$ rule. *Ann. Appl. Probab.* **5**(3) 809–833.
Van Mieghem, J. A. 2003. Due date scheduling: Asymptotic optimality of generalized longest queue and generalized largest delay rules. *Oper. Res.* **51**(1) 113–122.
Wallace, R. B., W. Whitt. 2005. A staffing algorithm for call centers with skill-based routing. *Manufacturing Service Oper. Management* **7**(4) 276–294.
Whitt, W. 2006. A multi-class fluid model for a contact center with skill-based routing. *Internat. J. Electronics Comm. (AEU)* **60**(2) 95–102.